

# Neuro Inspired Computing in FLASH

---



2014 Neuro Inspired Computational Elements Workshop

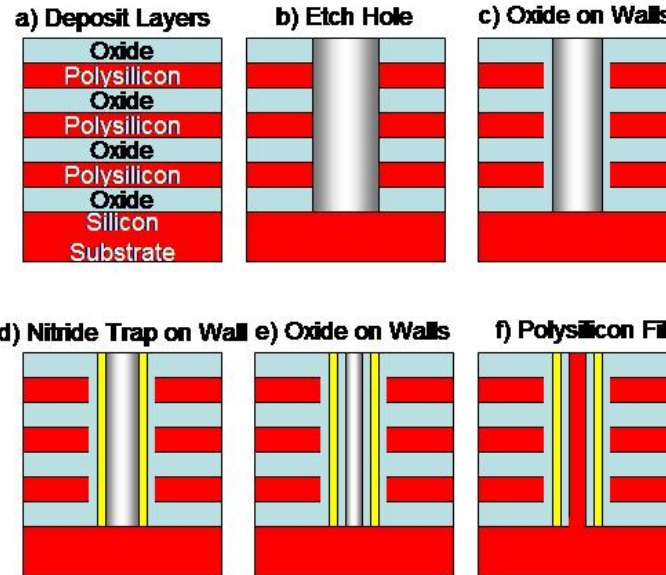
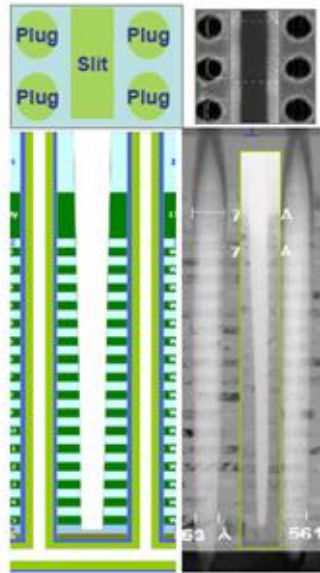
Kevin Gomez  
SSD Architecture

# Summary

- Embed specialized hardware into Solid State Drives (SSDs) for cortical processing assist
- Architecture modeling shows ~100x lower J/op compared to processing on host
- Identical to standard SSD - same manufacturing process and cost
- Re-purposed as cortical processor through firmware
- Open standards, e.g. same APIs as GPGPU, OpenCL, PyNN

# NAND Flash

has successfully transitioned to 3D



Jim Handy, 2013

An ingenious breakthrough which enables multiple layers of memory **without needing to pattern each layer**  
Decoupling NAND Flash production CAPEX from lithography (45% for planar down to 15% for 3D)\*

\*Samsung Analyst Day 2013 Memory Business

# 3D NAND Process

devices shipping in volume

# 3D NAND Process

devices shipping in volume

Silicon Substrate



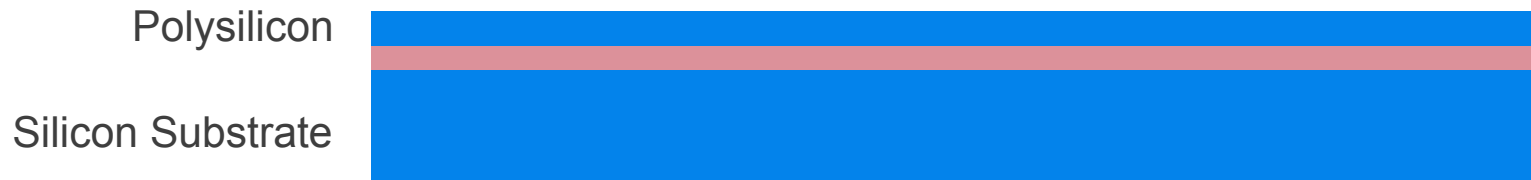
# 3D NAND Process

devices shipping in volume



# 3D NAND Process

devices shipping in volume



# 3D NAND Process

devices shipping in volume

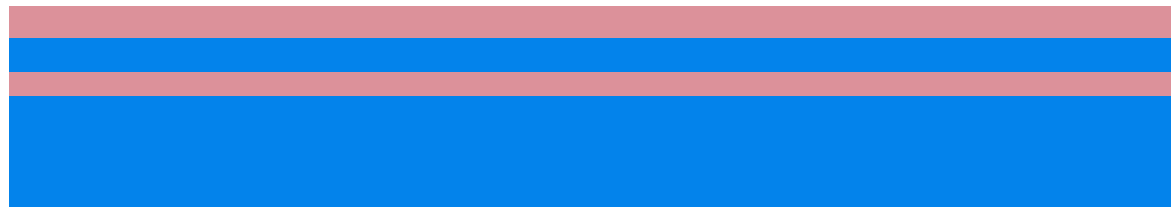




# 3D NAND Process

devices shipping in volume

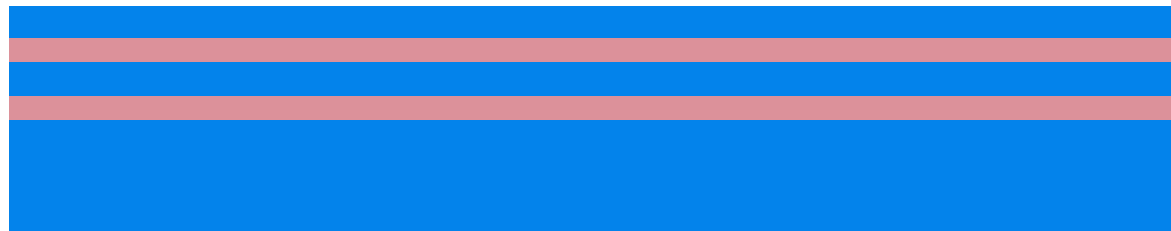
Silicon Substrate



# 3D NAND Process

devices shipping in volume

Silicon Substrate



# 3D NAND Process

devices shipping in volume

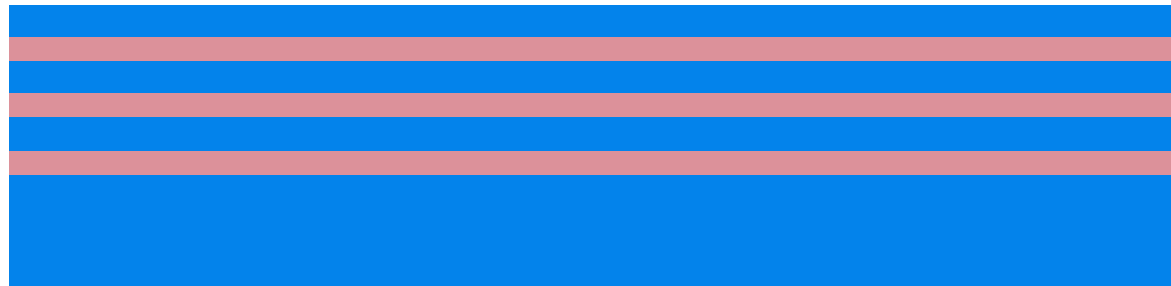
Silicon Substrate



# 3D NAND Process

devices shipping in volume

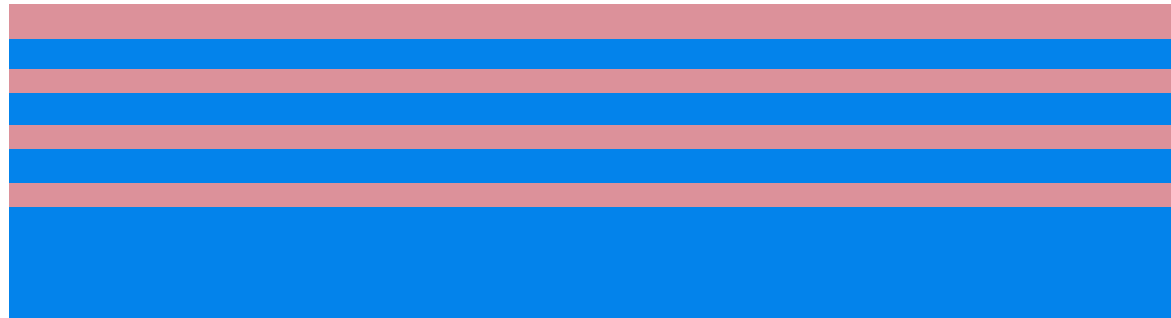
Silicon Substrate



# 3D NAND Process

devices shipping in volume

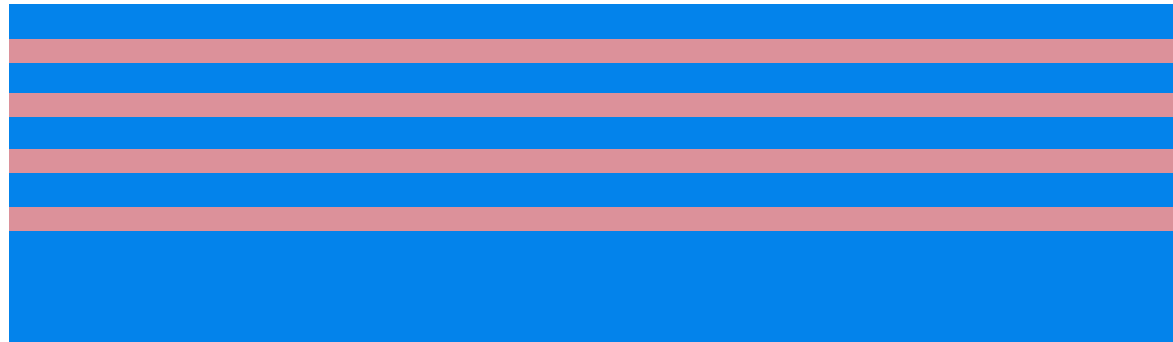
Silicon Substrate



# 3D NAND Process

devices shipping in volume

Silicon Substrate



# 3D NAND Process

devices shipping in volume

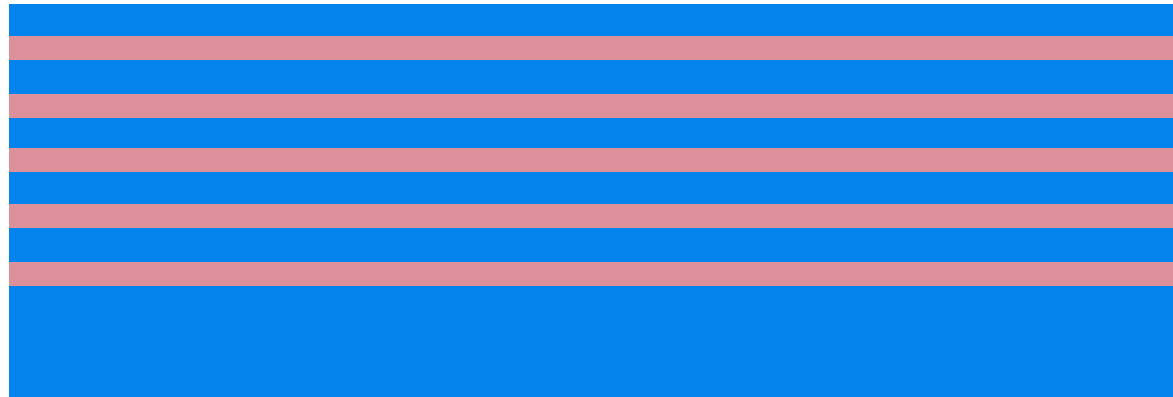
Silicon Substrate



# 3D NAND Process

devices shipping in volume

Silicon Substrate





# 3D NAND Process

devices shipping in volume

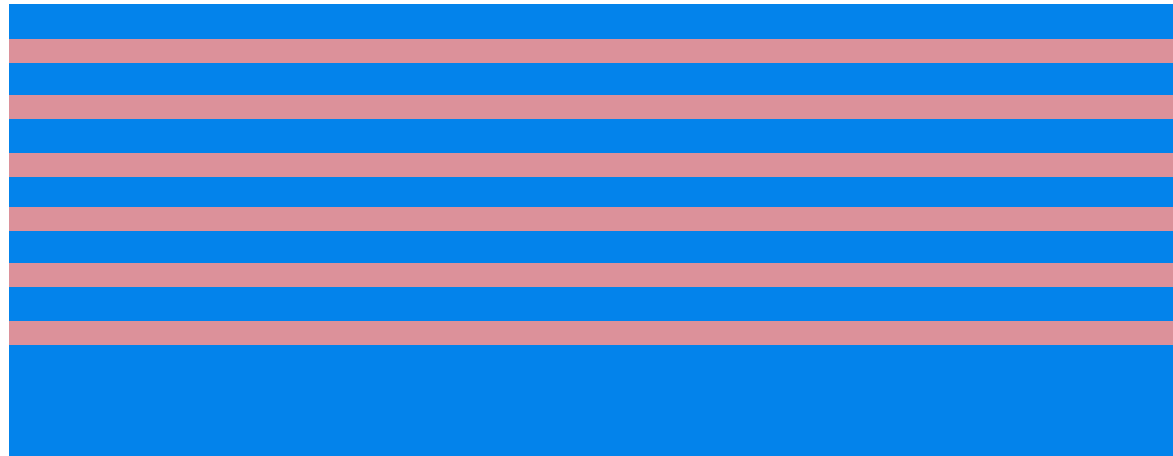
Silicon Substrate



# 3D NAND Process

devices shipping in volume

Silicon Substrate



# 3D NAND Process

devices shipping in volume

Silicon Substrate



# 3D NAND Process

devices shipping in volume

Silicon Substrate



# 3D NAND Process

devices shipping in volume

Silicon Substrate



# 3D NAND Process

devices shipping in volume

Silicon Substrate



# 3D NAND Process

devices shipping in volume

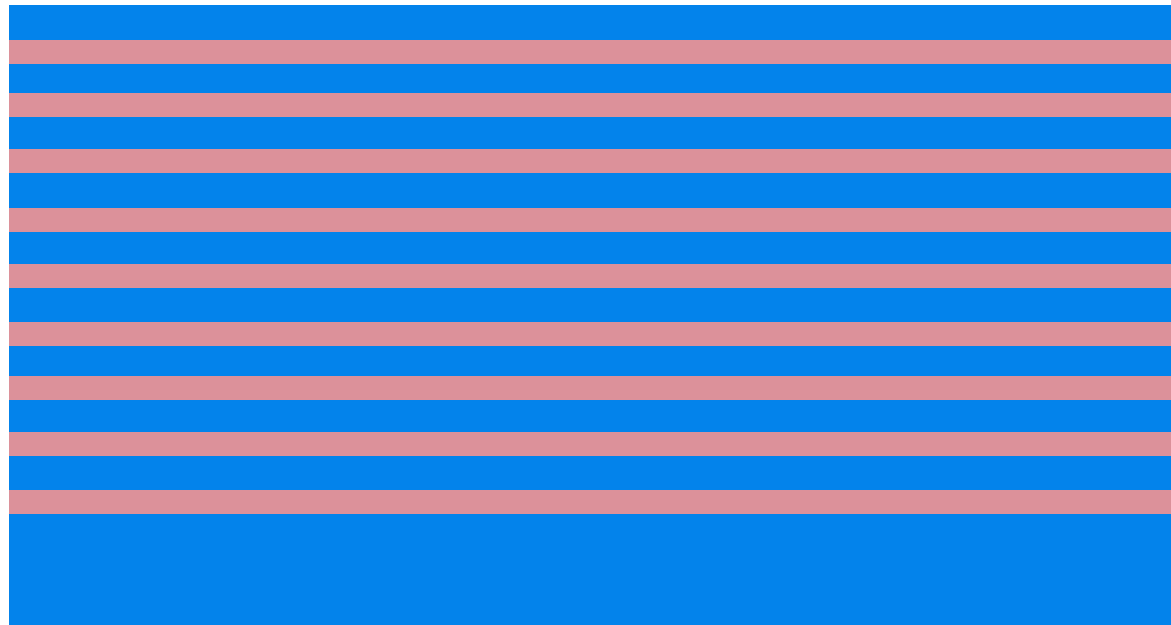
Silicon Substrate



# 3D NAND Process

devices shipping in volume

Silicon Substrate





# 3D NAND Process

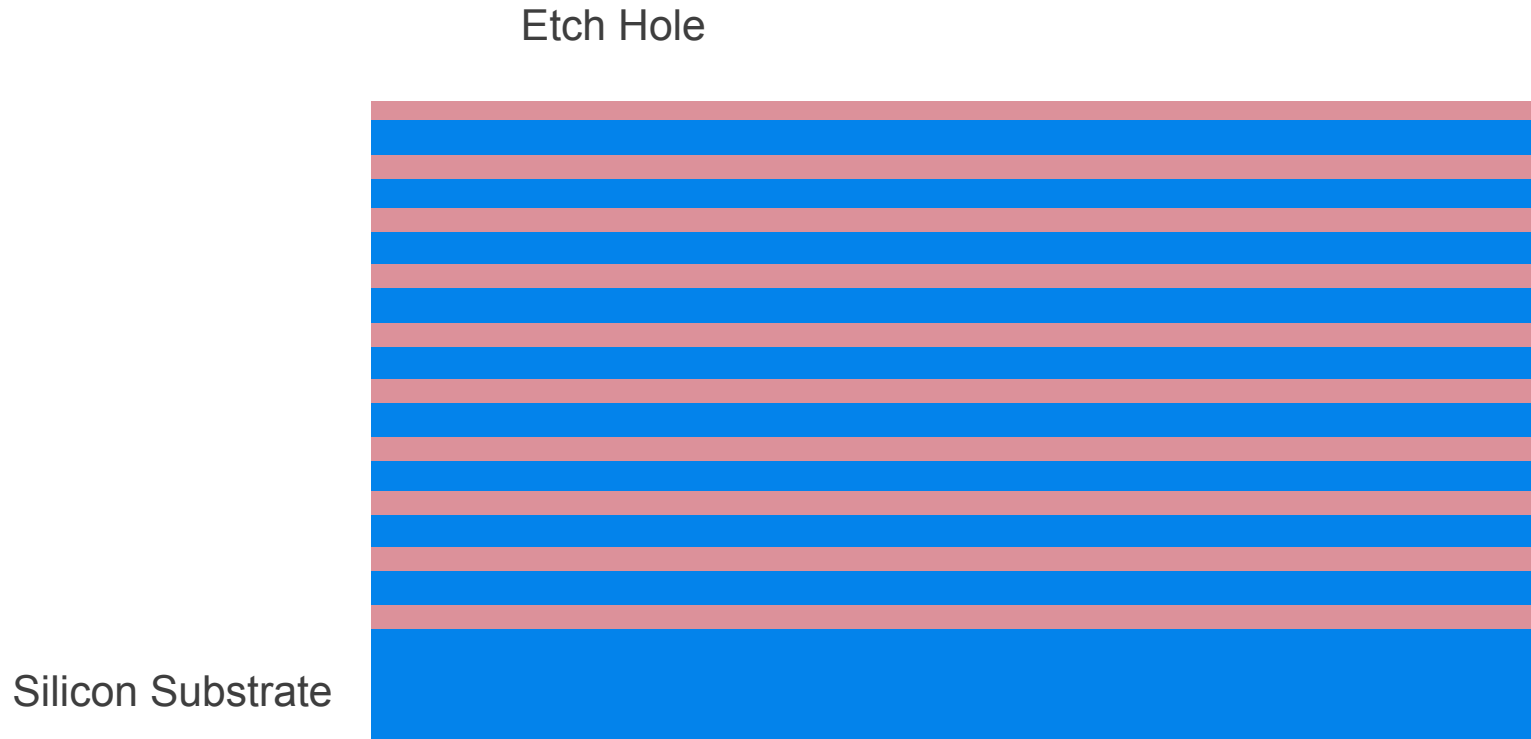
devices shipping in volume

Silicon Substrate



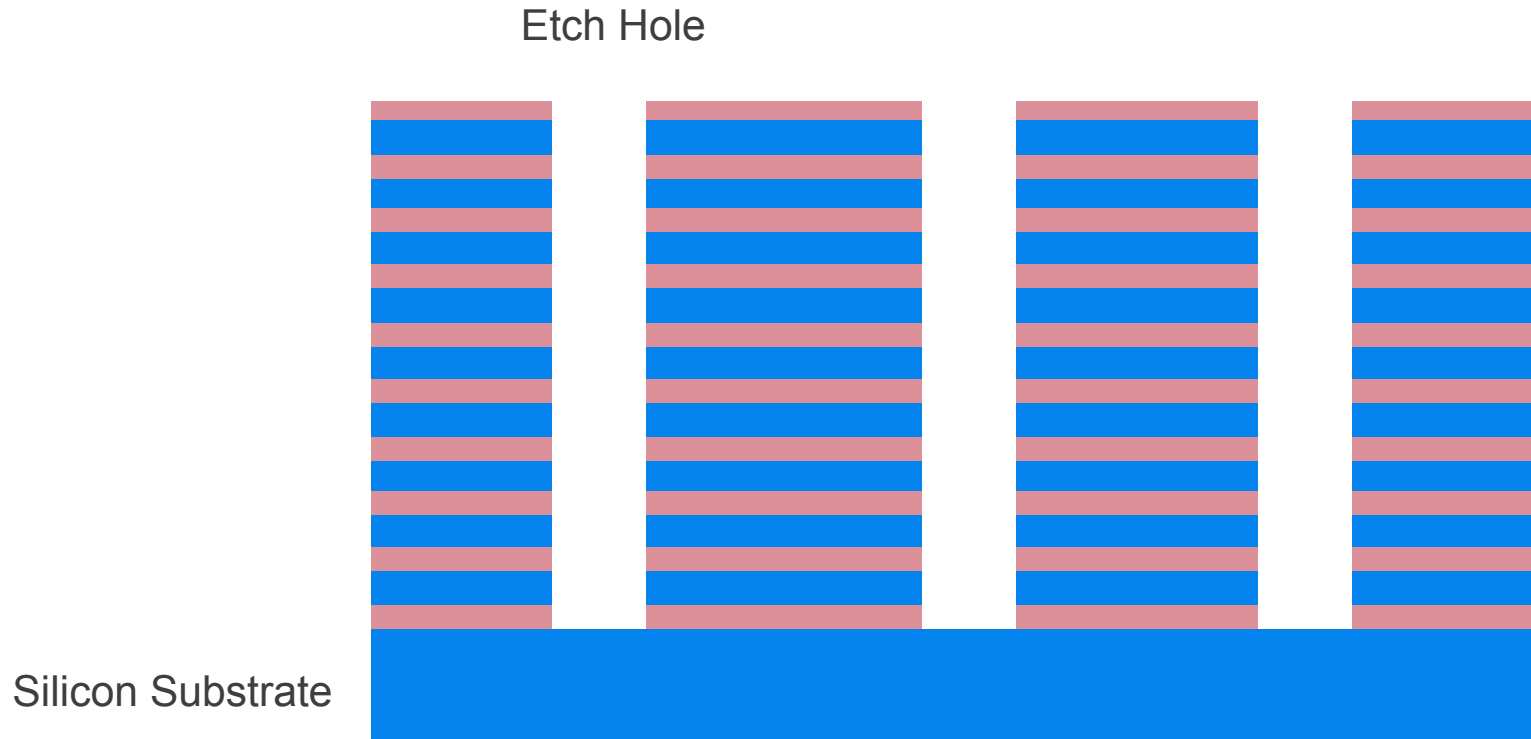
# 3D NAND Process

devices shipping in volume



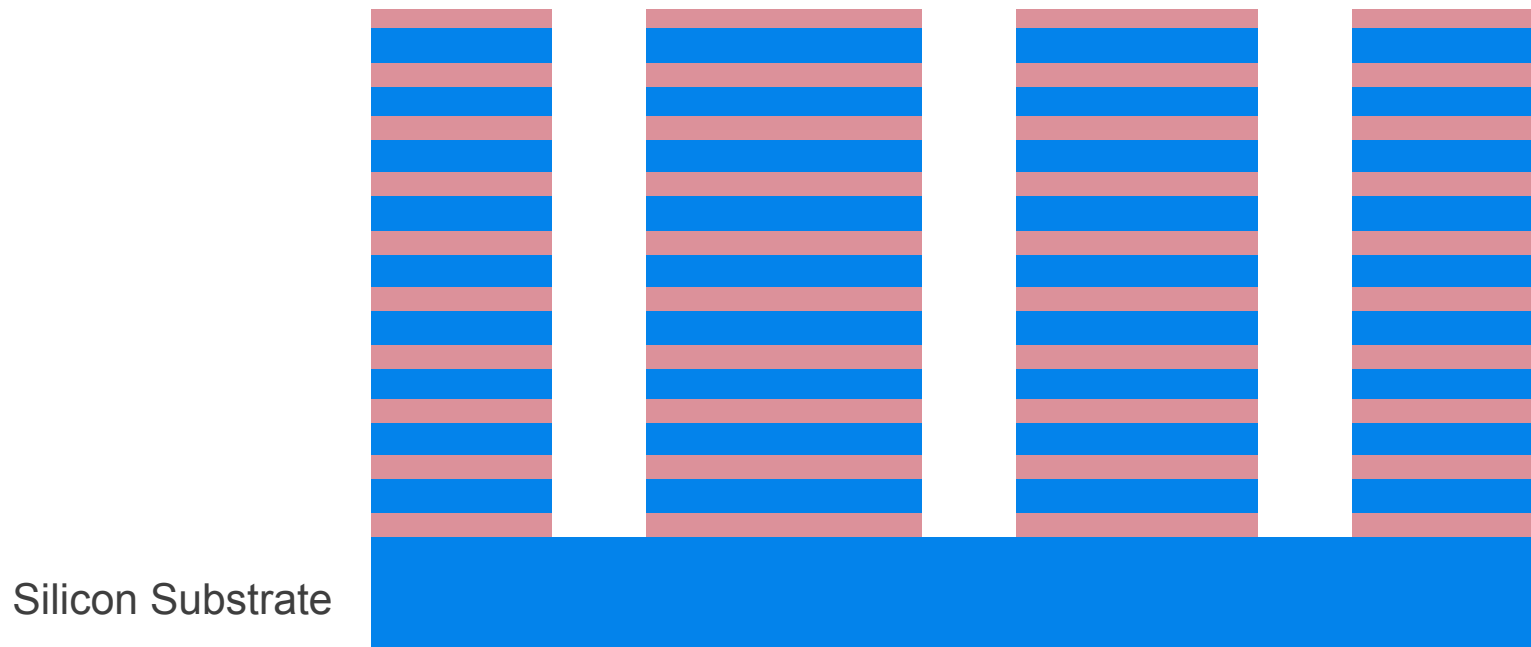
# 3D NAND Process

devices shipping in volume



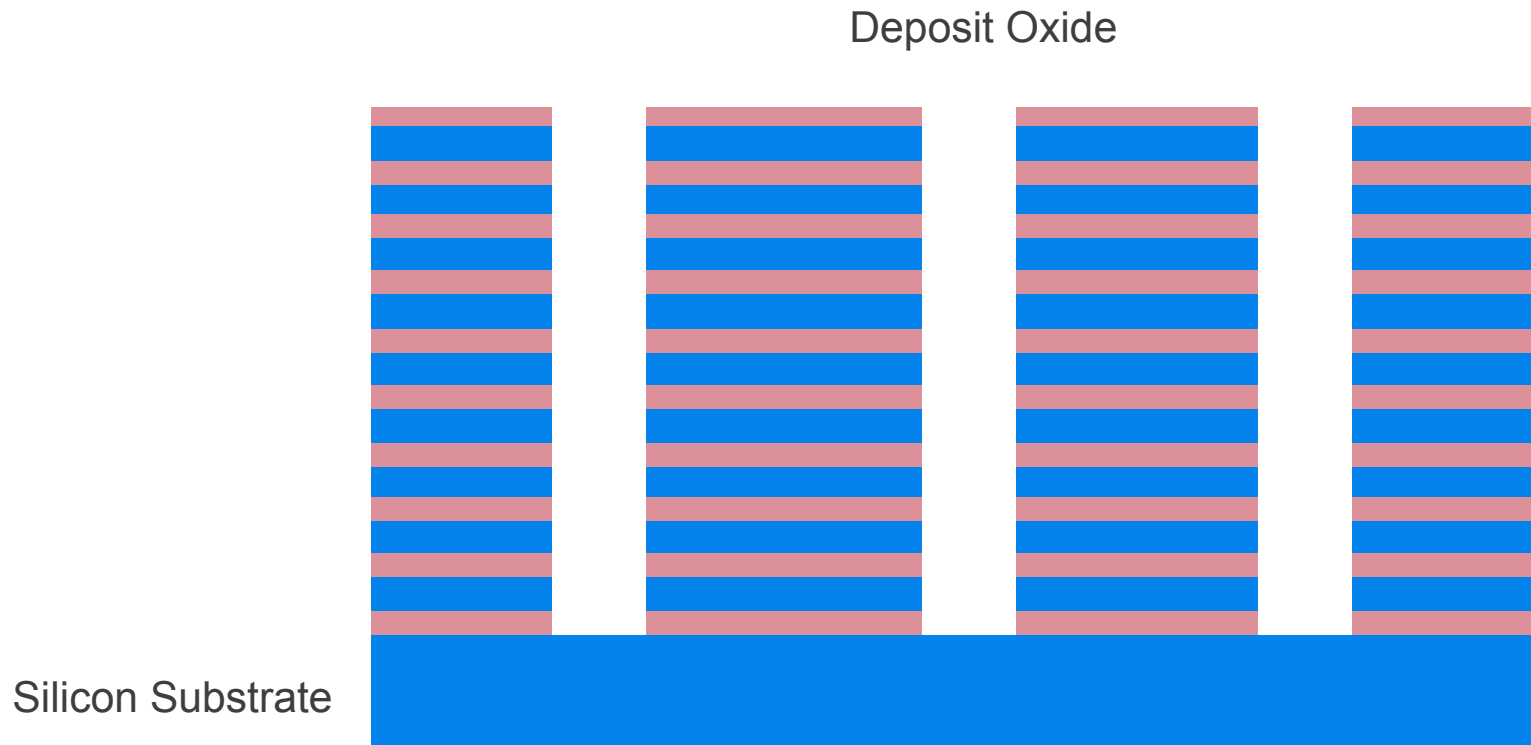
# 3D NAND Process

devices shipping in volume



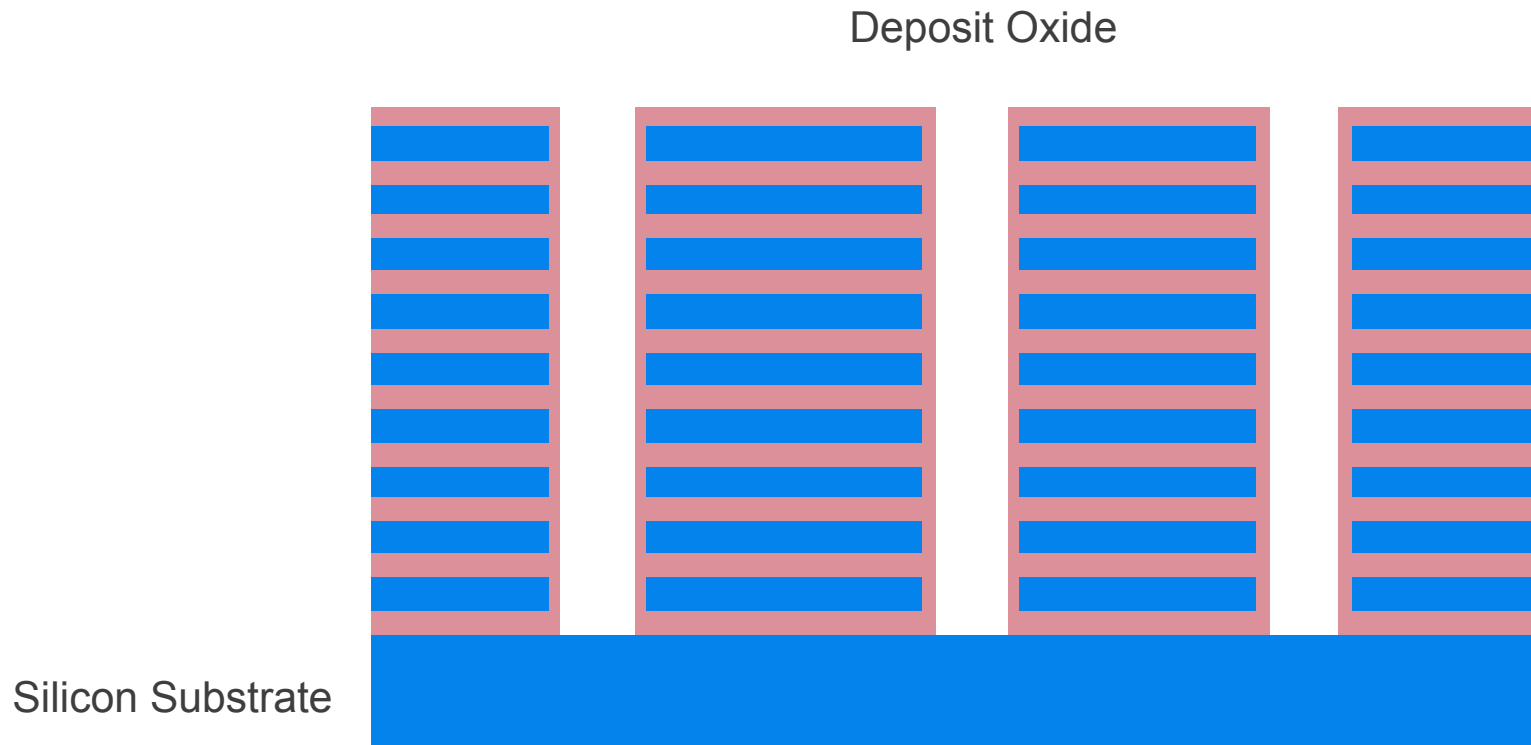
# 3D NAND Process

devices shipping in volume



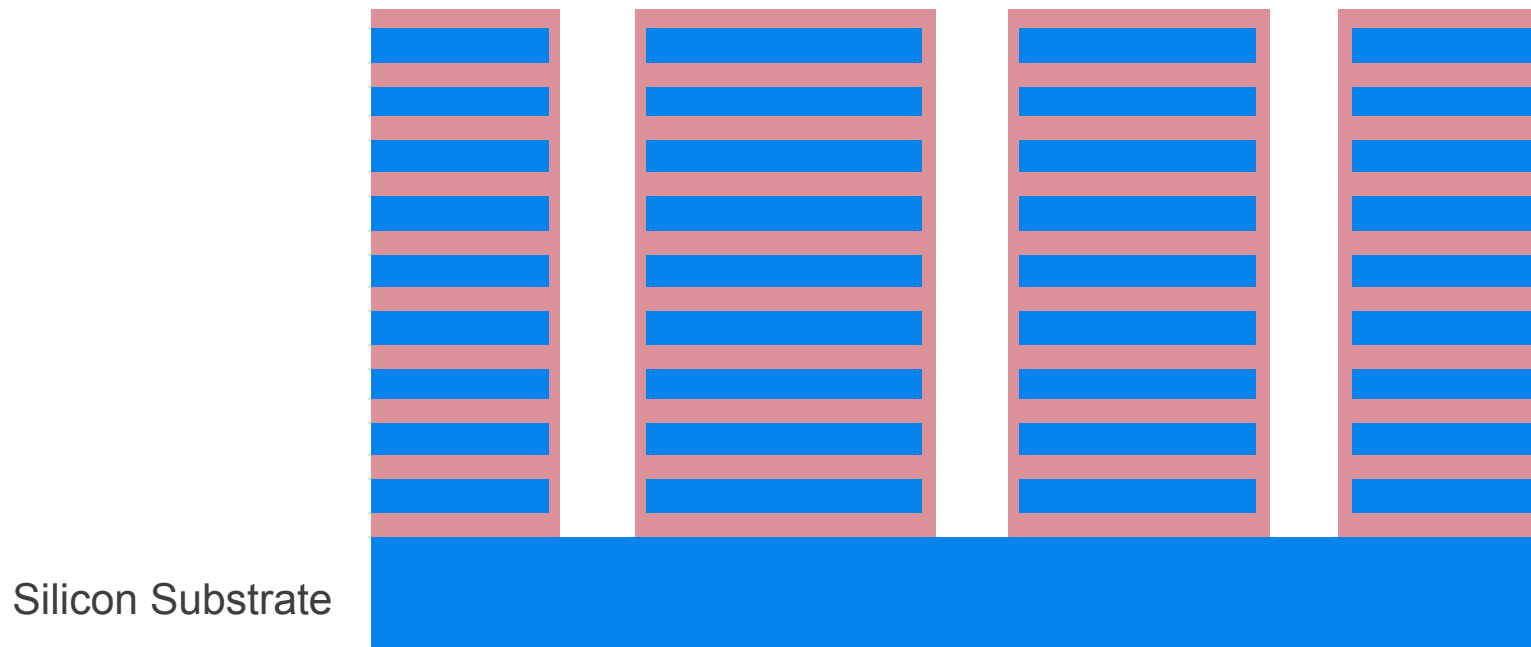
# 3D NAND Process

devices shipping in volume



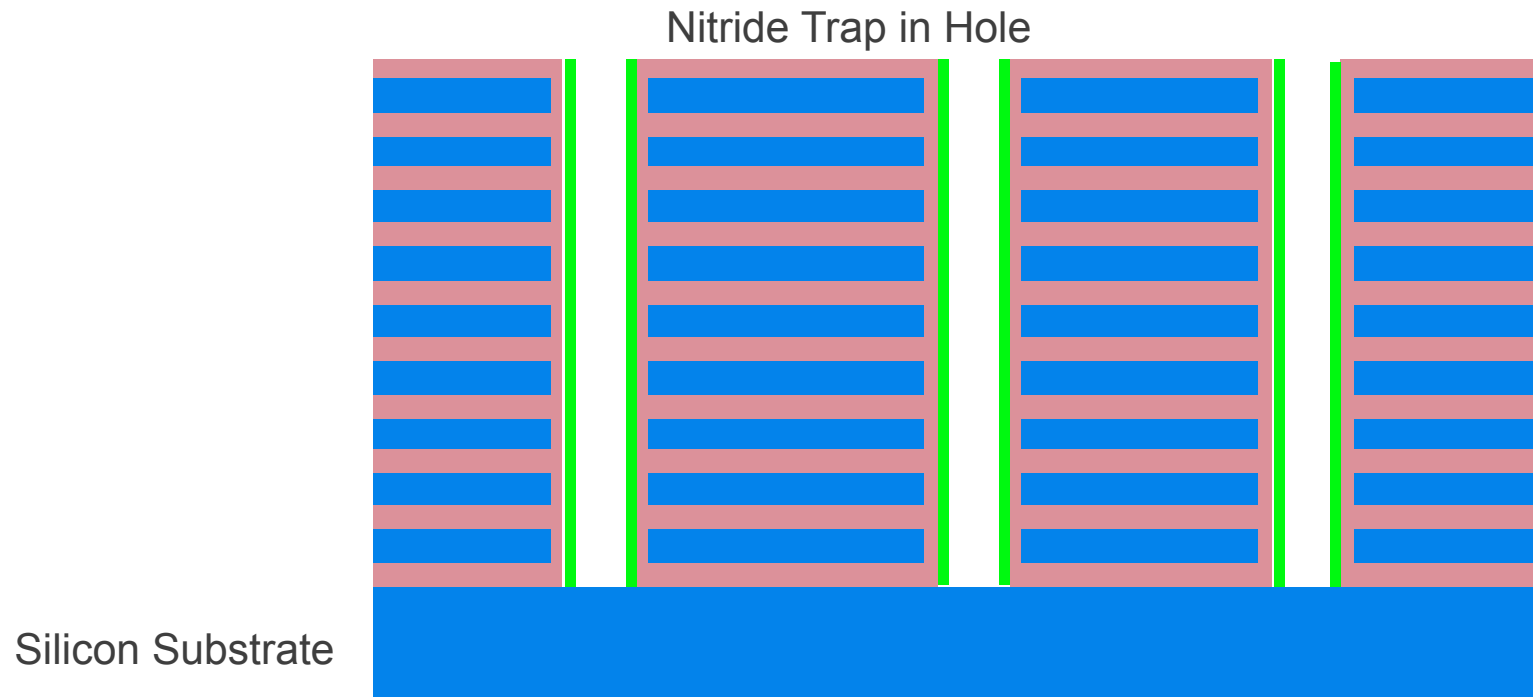
# 3D NAND Process

devices shipping in volume



# 3D NAND Process

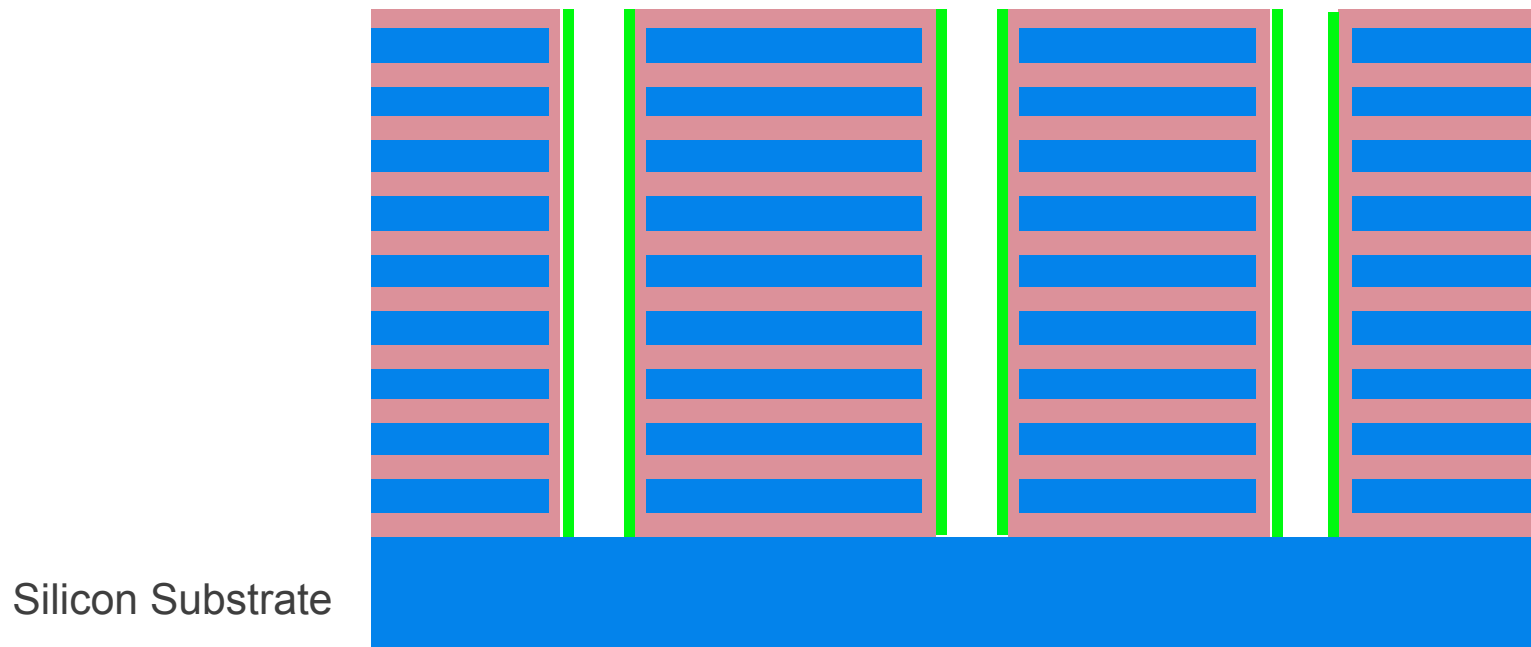
devices shipping in volume





# 3D NAND Process

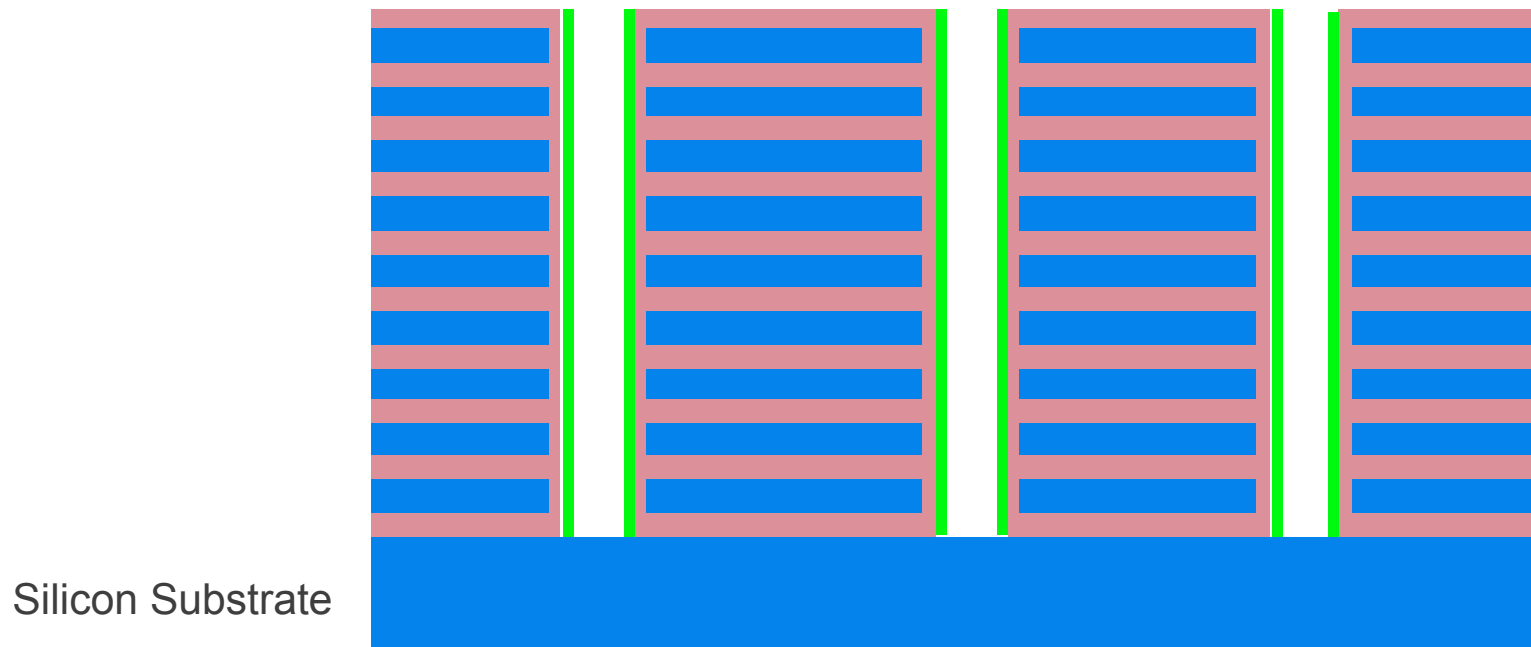
devices shipping in volume



# 3D NAND Process

devices shipping in volume

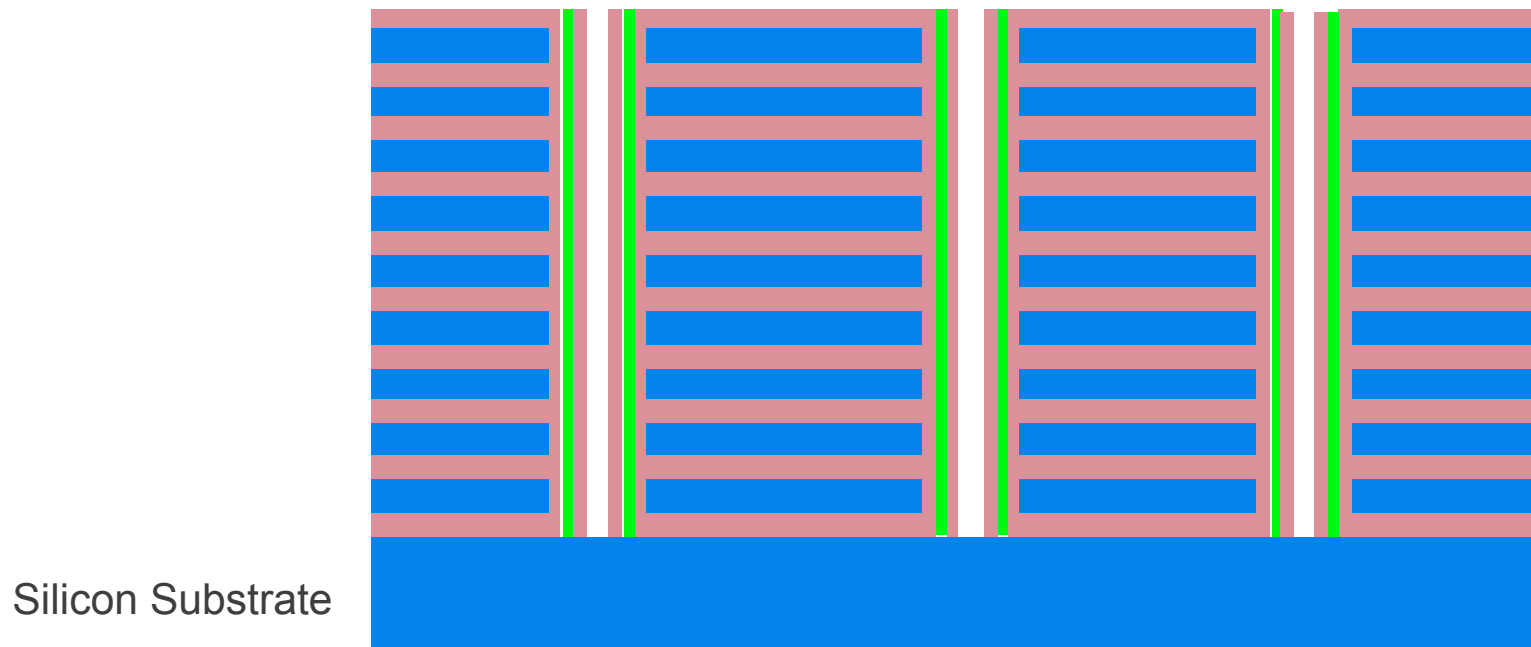
Deposit Oxide



# 3D NAND Process

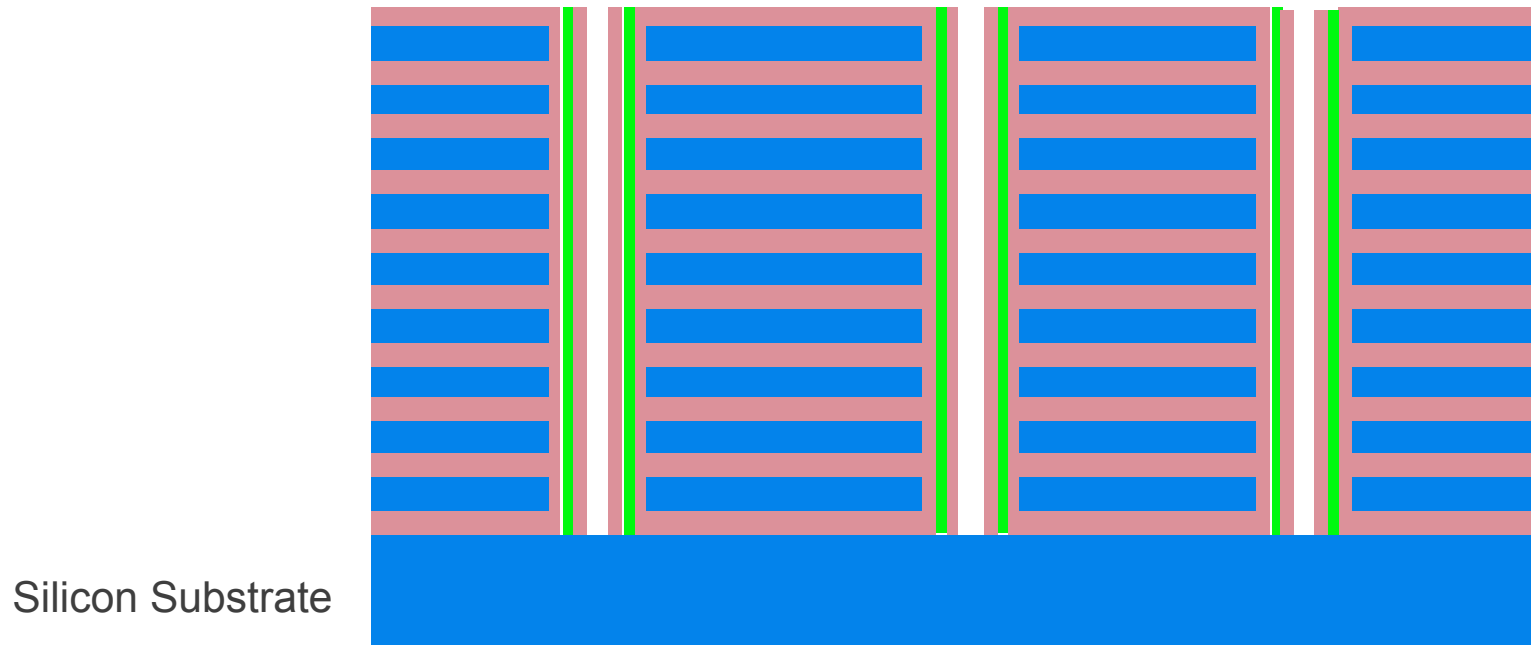
devices shipping in volume

Deposit Oxide



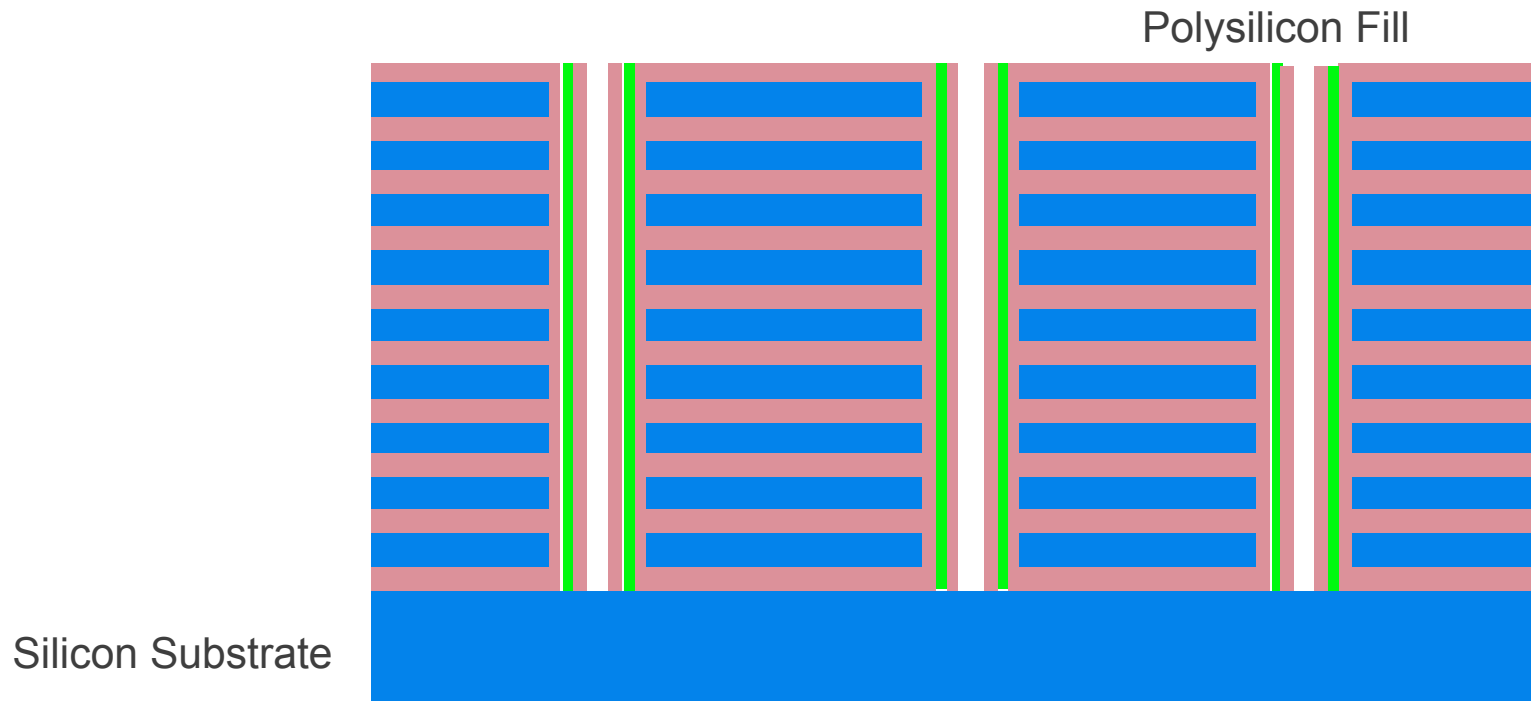
# 3D NAND Process

devices shipping in volume



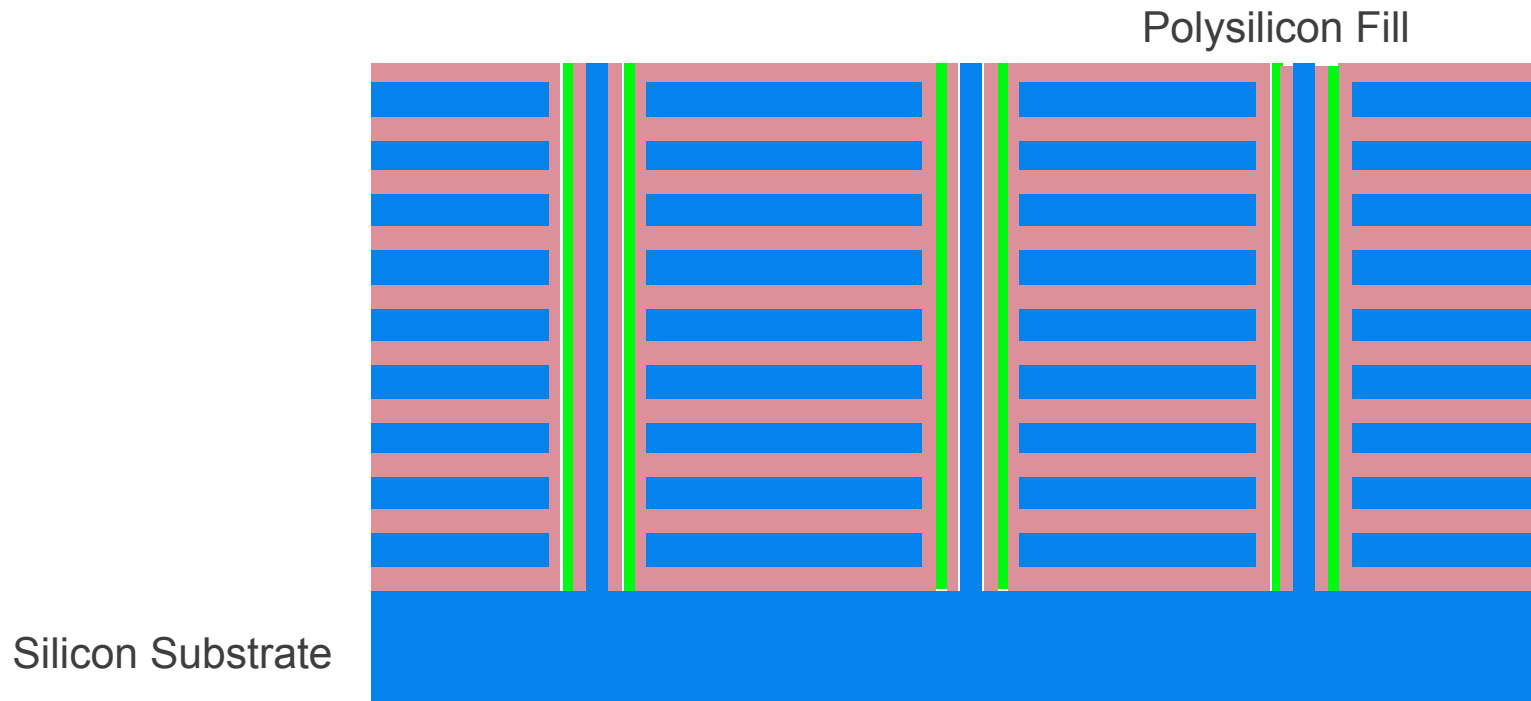
# 3D NAND Process

devices shipping in volume



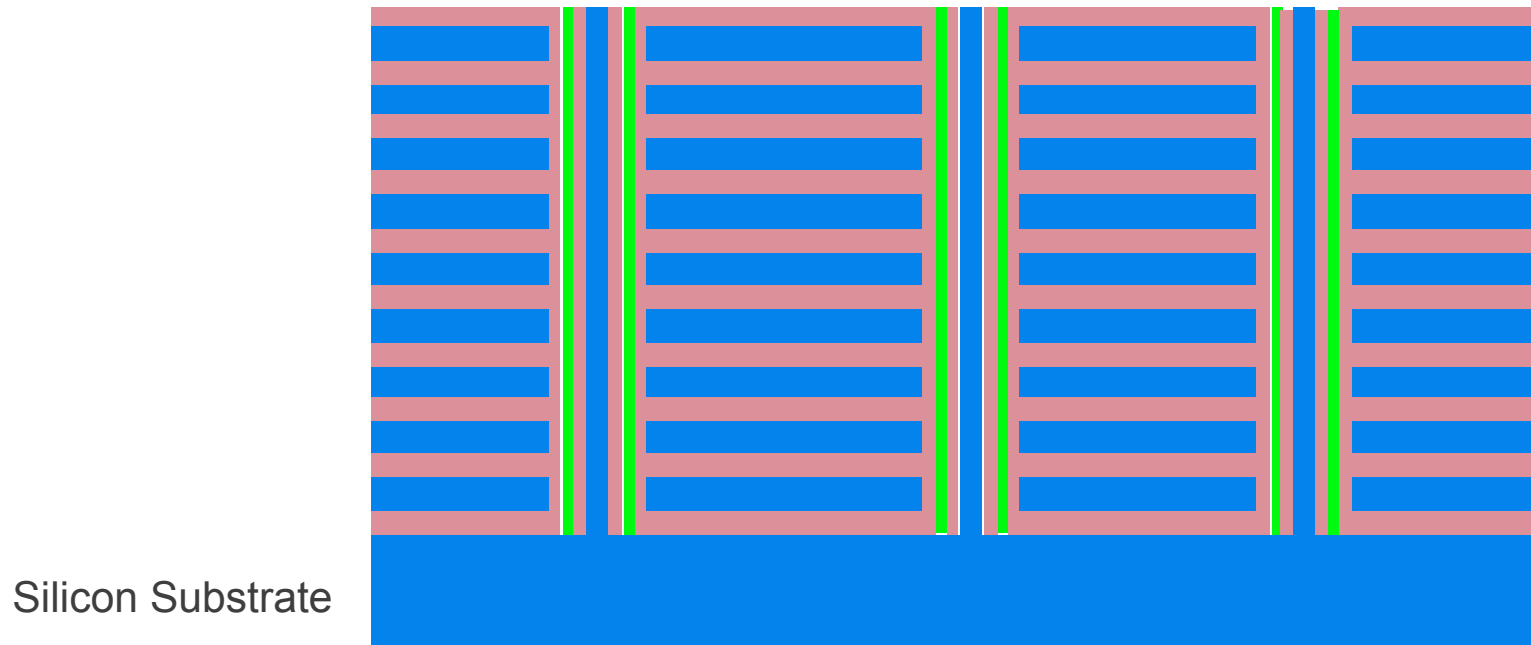
# 3D NAND Process

devices shipping in volume



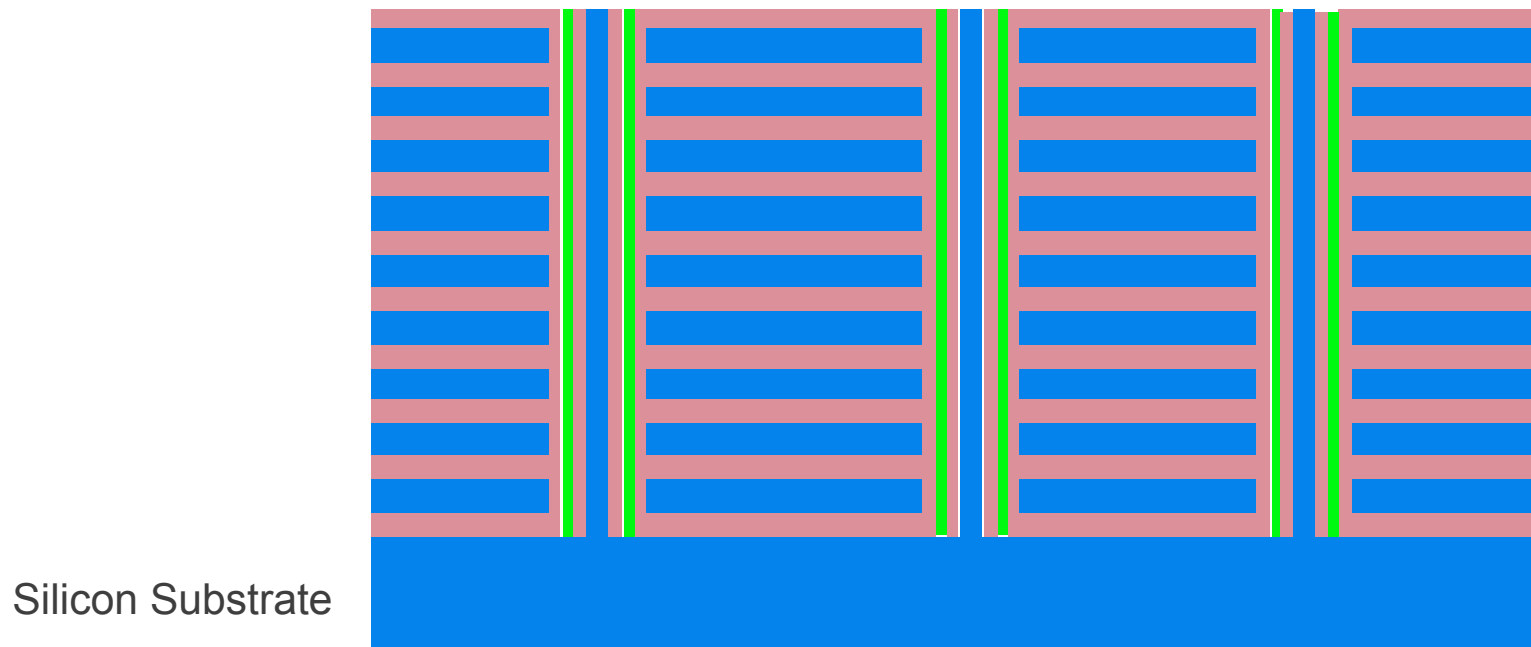
# 3D NAND Process

devices shipping in volume



# 3D NAND Process

devices shipping in volume



**“The burden will shift from lithography to deposition and etch”**  
- Ritu Shrivastava, Sandisk



# ITRS – Technology Trends

for DRAM and FLASH Memory

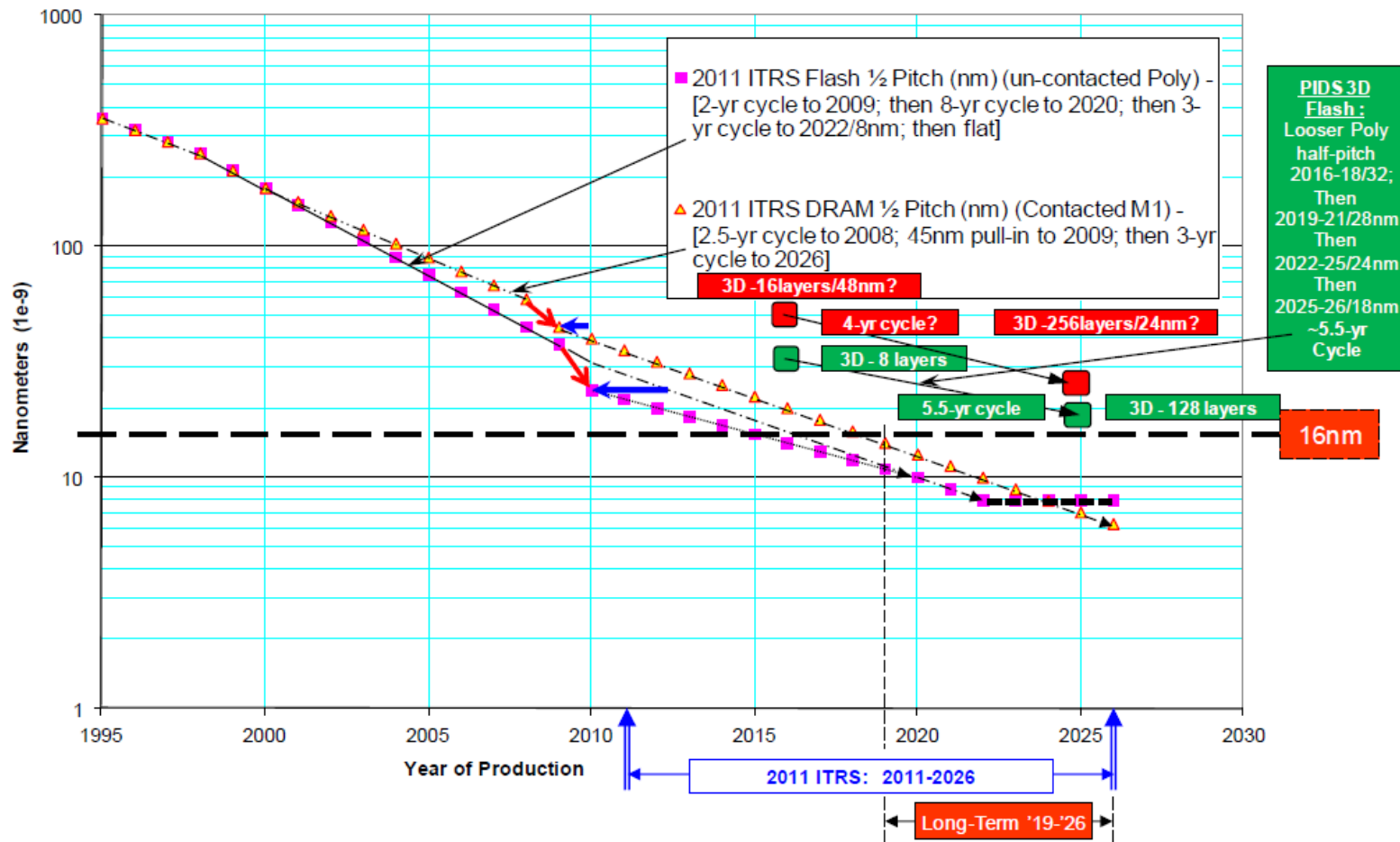


Figure 10 2011 ITRS—DRAM and Flash Memory Half Pitch Trends

# NAND Flash Scaling - ITRS

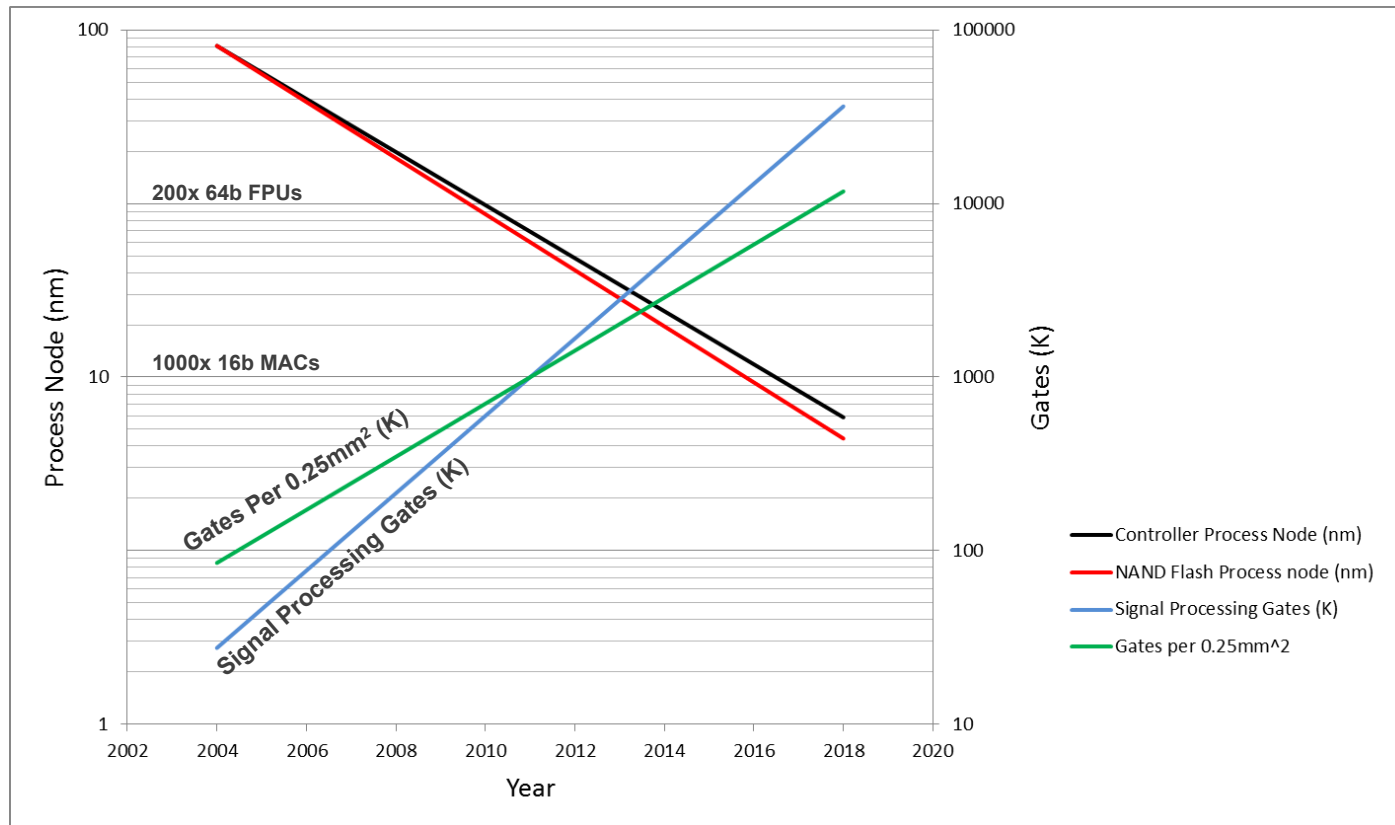
Relax planar scaling, push into 3rd dimension, continue Moore's law

| <b>NAND Flash</b>  |             |             |              |              |                    |                    |                    |                  |                  |                  |                |                |                |                |
|--|-------------|-------------|--------------|--------------|--------------------|--------------------|--------------------|------------------|------------------|------------------|----------------|----------------|----------------|----------------|
| <i>Year of Production</i>  | <i>2012</i> | <i>2013</i> | <i>2014</i>  | <i>2015</i>  | <i>2016</i>        | <i>2017</i>        | <i>2018</i>        | <i>2019</i>      | <i>2020</i>      | <i>2021</i>      | <i>2022</i>    | <i>2023</i>    | <i>2024</i>    | <i>2025</i>    |
| <i>Uncontacted poly 1/2 pitch (nm)</i>                               | <b>20</b>   | <b>18</b>   | <b>17</b>    | <b>15</b>    | <b>14</b>          | <b>13</b>          | <b>12</b>          | <b>11</b>        | <b>10</b>        | <b>9</b>         | <b>8</b>       | <b>8</b>       | <b>8</b>       | <b>8</b>       |
| <i>Number of word lines in one NAND string</i>                       | <b>64</b>   | <b>64</b>   | <b>64</b>    | <b>64</b>    | <b>64</b>          | <b>64</b>          | <b>64</b>          | <b>64</b>        | <b>64</b>        | <b>64</b>        | <b>64</b>      | <b>64</b>      | <b>64</b>      | <b>64</b>      |
| <i>Dominant Cell type</i>  | <b>FG</b>   | <b>FG</b>   | <b>FG/CT</b> | <b>FG/CT</b> | <b>CT-3D</b>       | <b>CT-3D</b>       | <b>CT-3D</b>       | <b>CT-3D</b>     | <b>CT-3D</b>     | <b>CT-3D</b>     | <b>CT-3D</b>   | <b>CT-3D</b>   | <b>CT-3D</b>   | <b>CT-3D</b>   |
| <i>Maximum number of bits per chip (SLC/MLC)</i>                     |             |             |              |              | <b>128G / 256G</b> | <b>256G / 512G</b> | <b>256G / 512G</b> | <b>512G / 1T</b> | <b>512G / 1T</b> | <b>512G / 1T</b> | <b>1T / 2T</b> | <b>1T / 2T</b> | <b>1T / 2T</b> | <b>2T / 4T</b> |
| <i>Minimum array 1/2 pitch - F(nm) [15]</i>                          |             |             |              |              | <b>32nm</b>        | <b>32nm</b>        | <b>32nm</b>        | <b>28nm</b>      | <b>28nm</b>      | <b>28nm</b>      | <b>24nm</b>    | <b>24nm</b>    | <b>24nm</b>    | <b>18nm</b>    |
| <i>Number of 3D layers for array at minimum 1/2 array pitch [16]</i> |             |             |              |              | <b>8</b>           | <b>16</b>          | <b>32</b>          | <b>32</b>        | <b>64</b>        | <b>64</b>        | <b>98</b>      | <b>98</b>      | <b>98</b>      | <b>128</b>     |

- ITRS Winter Public Conference Dec 2012 Hsinchu, Taiwan

# Sustaining NAND density growth

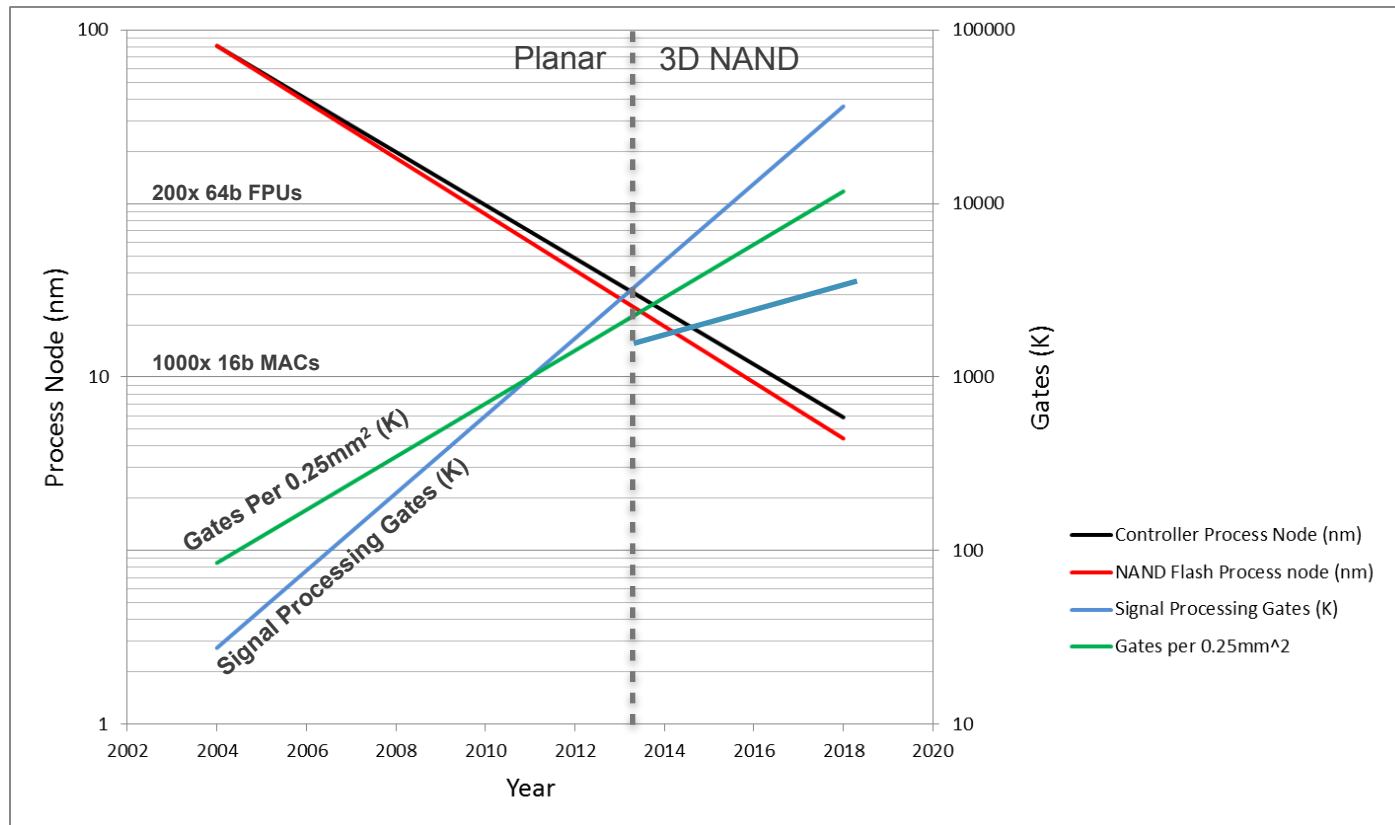
Near Shannon limit, Iterative Low-Density Parity-Check channels >1M gates



Li, Peng, Kevin Gomez, and David J. Lilja. "Exploiting Free Silicon for Energy-Efficient Computing Directly in NAND Flash-based Solid-State Storage Systems." *IEEE HPEC 2013*

# Sustaining NAND density growth

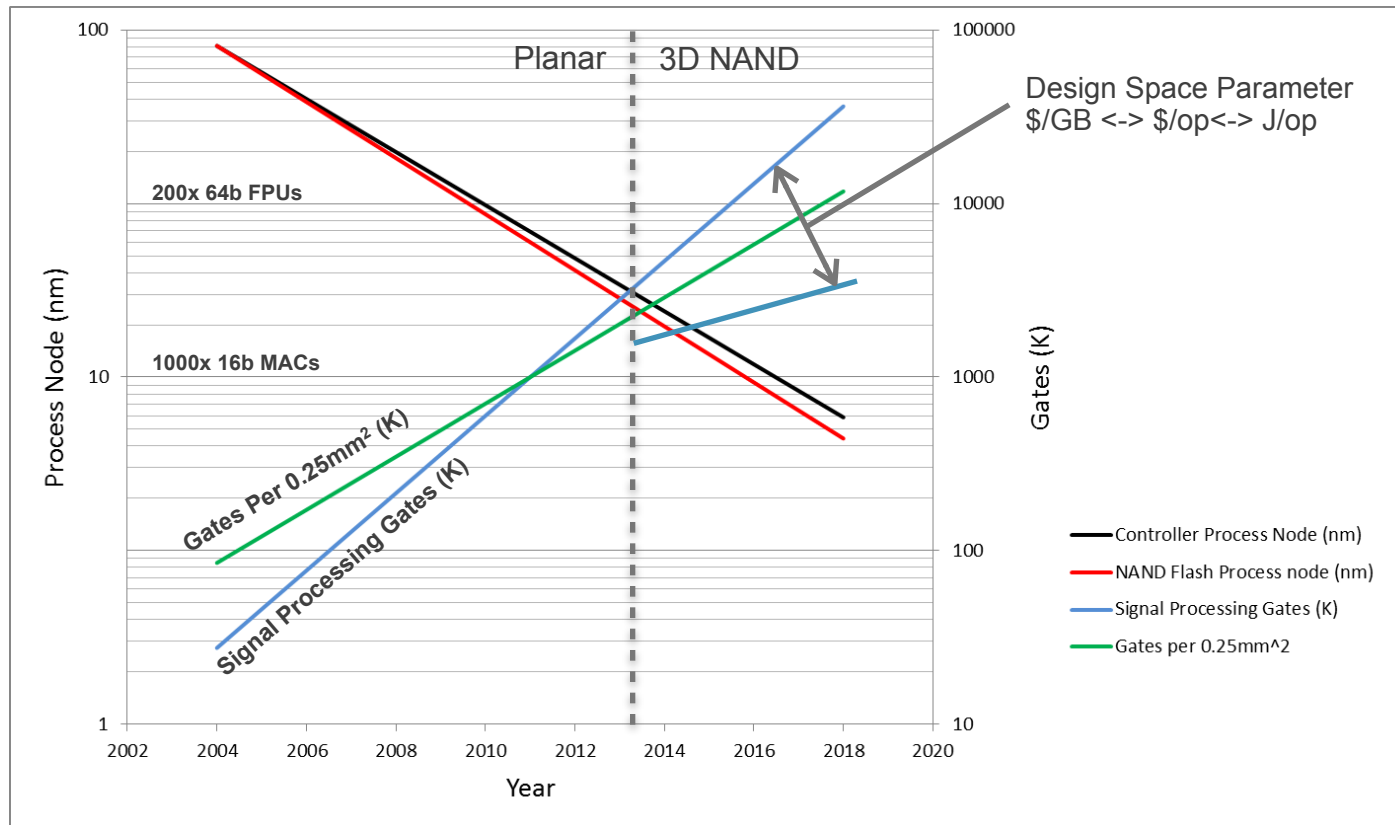
Near Shannon limit, Iterative Low-Density Parity-Check channels >1M gates



Li, Peng, Kevin Gomez, and David J. Lilja. "Exploiting Free Silicon for Energy-Efficient Computing Directly in NAND Flash-based Solid-State Storage Systems." IEEE HPEC 2013

# Sustaining NAND density growth

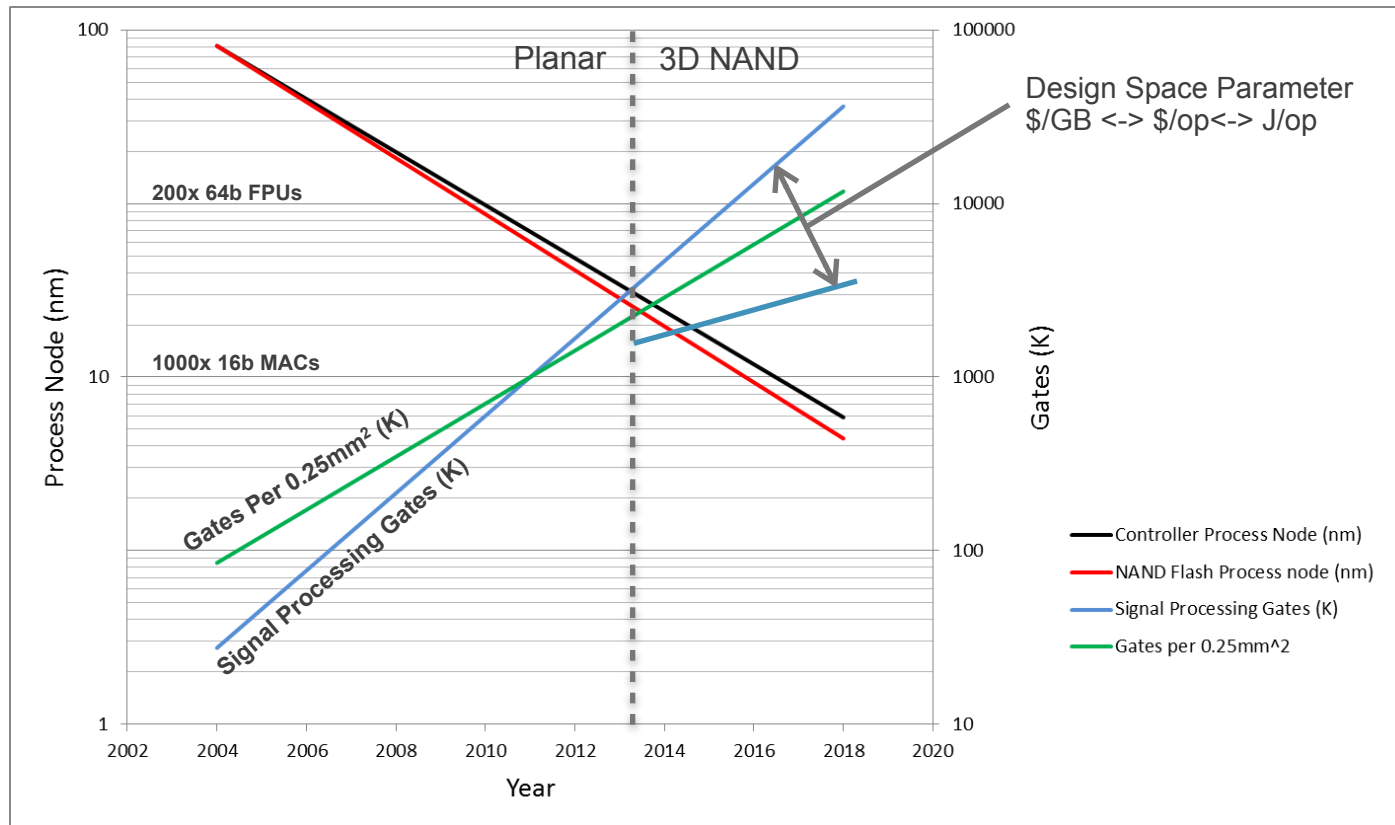
Near Shannon limit, Iterative Low-Density Parity-Check channels >1M gates



Li, Peng, Kevin Gomez, and David J. Lilja. "Exploiting Free Silicon for Energy-Efficient Computing Directly in NAND Flash-based Solid-State Storage Systems." *IEEE HPEC* 2013

# Sustaining NAND density growth

Near Shannon limit, Iterative Low-Density Parity-Check channels >1M gates

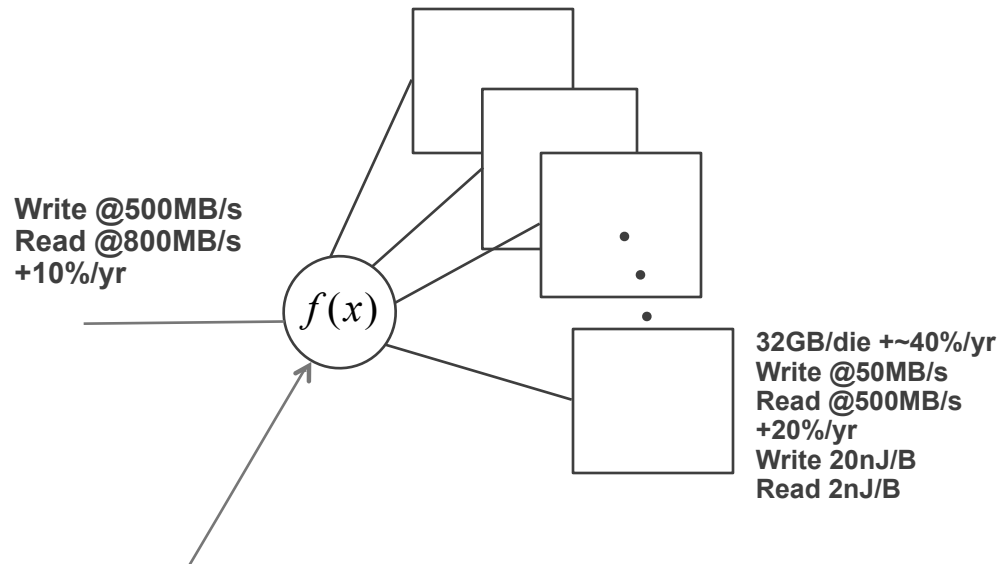


Li, Peng, Kevin Gomez, and David J. Lilja. "Exploiting Free Silicon for Energy-Efficient Computing Directly in NAND Flash-based Solid-State Storage Systems." *IEEE HPEC 2013*

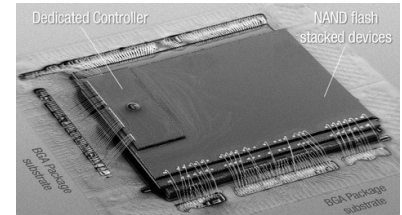
**Cost of adding specialized cortical hardware automation is marginal**

# Scaling up a Cortical processor

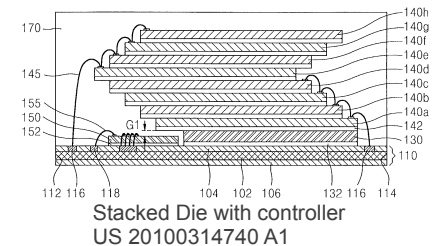
4 to 16 Flash Die per Package



Currently ~2M gates for LDPC and packet switching  
Substrate for specialized or reconfigurable cortical hardware



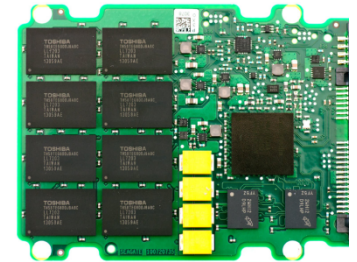
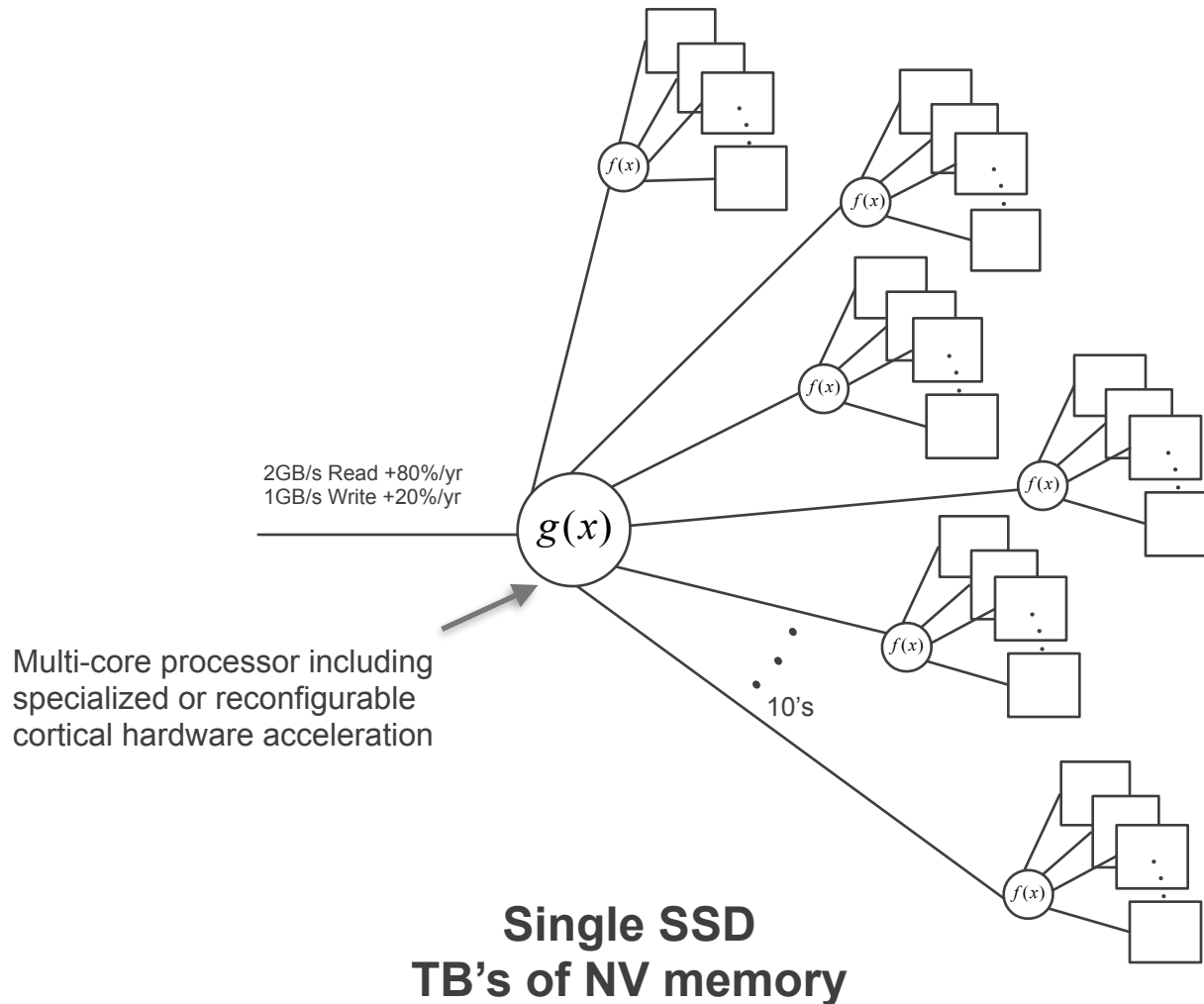
Example eMMC device (Micron)



**Single NAND Flash package ~ 5TB/in<sup>3</sup> +40%/yr**

# Scaling up a Cortical processor

10's of Flash Packages in each SSD

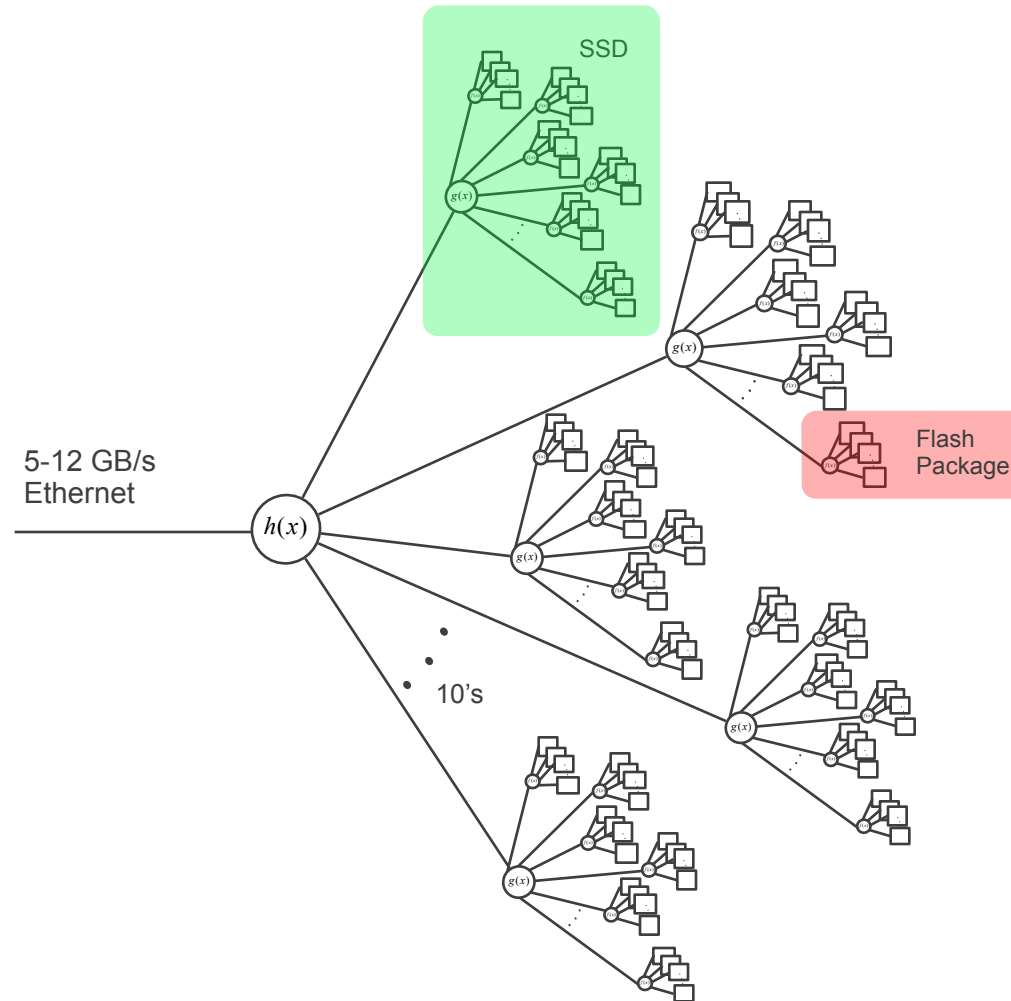


Seagate 600 Pro SSD



# Scaling up a Cortical processor

Array of SSDs in 1U Rack in Cloud Compute Server



# Convergence of Trends

## Why Now

- Flash Memory cost and SNR driven by mobile computing market
  - Increased investment in signal processing silicon at NAND interface – low marginal cost for added compute
- Power Wall – end of Dennard scaling (power  $\propto 1/L$  instead of  $1/L^3$ )
  - Since 2005 has driven multi-core parallelism to maintain compute cost-performance trajectory
  - In turn has forced parallel programming into the mainstream
- Moore's Law post Power Wall continues to provide gates at  $1/L^2$  which can not all be switching simultaneously
  - Increased adoption of power islanded heterogeneous architectures operating at device power budget
- Memory Wall – exponentially growing gap between processor and memory performance
  - Continues to drive tighter integration of memory and compute. GPU processing is a temporary reprieve

# Heterogeneous Architectures

Lots of efficient H/W automation – powered off most of the time

As Moore continues to increase the number of transistors on silicon at a scale of  $1/L^2$  while power is only decreasing as  $1/L$  ...

... we can afford to 'overprovision' the chip – i.e. use the TDP (total die power budget) using just a subset of the chip's resources – for example use the entire budget on compute while shutting down global on-chip communication resources.

Enables peak performance (using all available power) on diverse workloads.

This may signal that the right time for Reconfigurable Computing has arrived – specialized hardware acceleration, powered off most of the time.

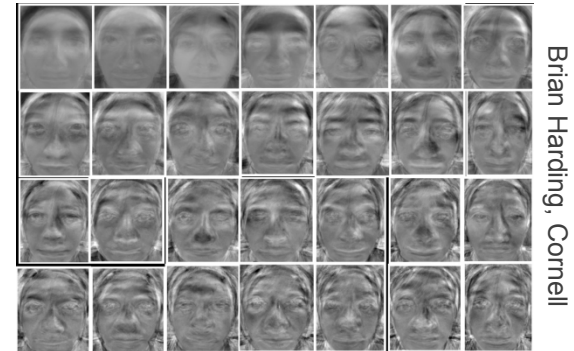
# Why NAND Flash and not other NVM technologies

- NAND is a block device and requires a significant and growing investment in signal processing to enable it's continued scaling
- This signal processing overhead is best situated close to NAND to minimize the energy cost of data movement
- NAND has no delusions of being a DRAM replacement like PCM or STTRAM with low-latency and close to byte addressable architectures which will not tolerate any significant signal processing overhead
- **It is not about the technology – it's the economics** - SSDs exist due to the demand for consumer grade NAND devices for the smartphone, tablet, SD Card and USB memory markets.
- Cortical Inspired Compute Elements embedded in SSDs likewise will succeed or fail purely on economics (\$/op, J/op) not technology

# Architecture Modeling

Facial recognition task which is a proxy algorithm for content based image retrieval:

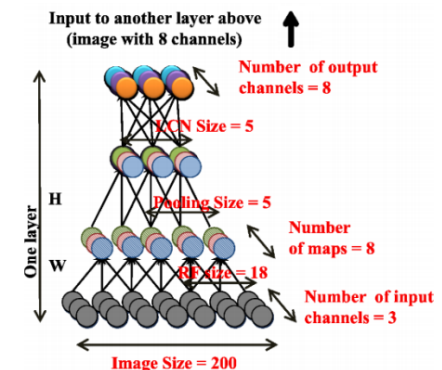
Compute on 16 channel SSD is  $\sim 0.2\text{mJ}/\text{face}$   
**150X lower J/face** than computing on Host



SCE stores and computes on eigenfaces

Boltzmann machine task- a proxy for many machine learning and data intensive scientific compute algorithms:

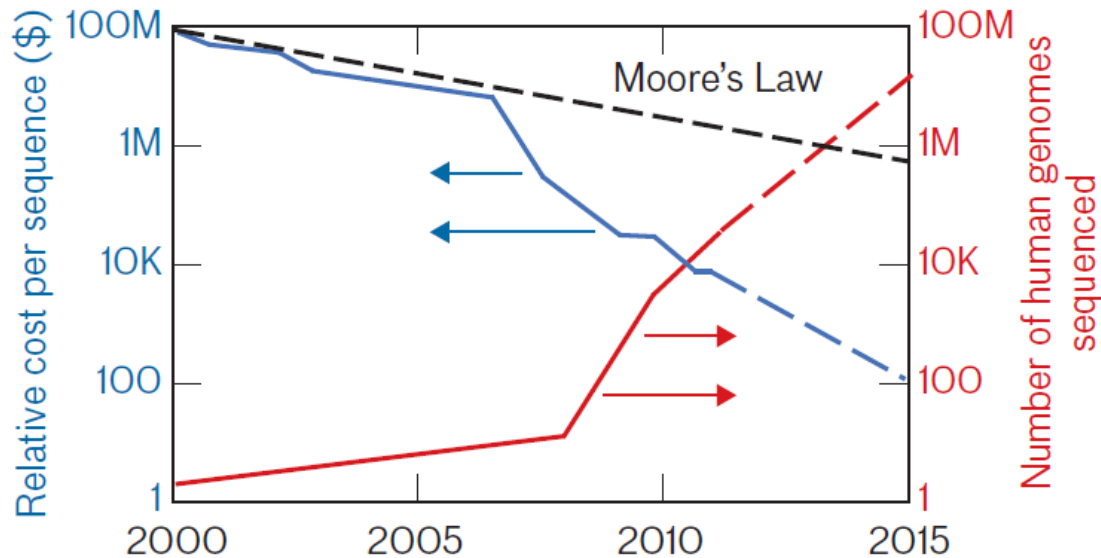
Compute on SSD is  $\sim 40\text{X}$  lower J/Op compared to Quad-Core host



“Building High-level Features Using Large Scale Unsupervised Learning”,  
Quok Le et al, 2012

# The need for Energy Efficiency

Big Data Analytics is no longer a Niche



Advances in DNA sequencing are rapidly decreasing the cost of whole human genome sequencing

As a result, the number of humans being sequenced is increasing significantly

Data needing to be processed is rapidly outpacing computing performance-cost.

Together these drive the need for greater efficiency.

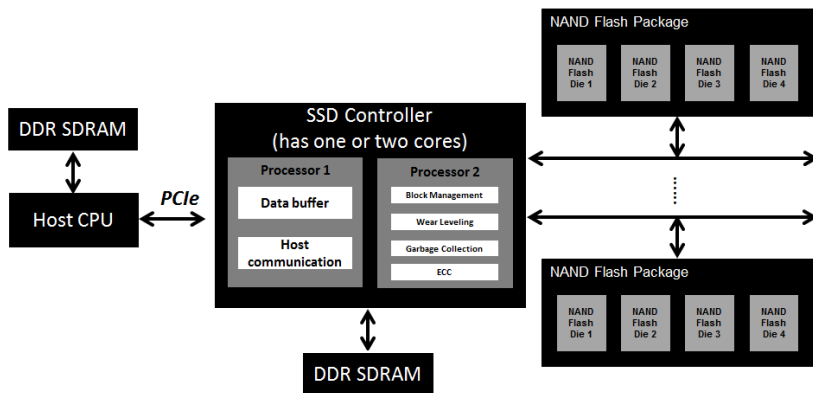
*"Taming Biological Data with D4M", Kepner 2013*

Thank You

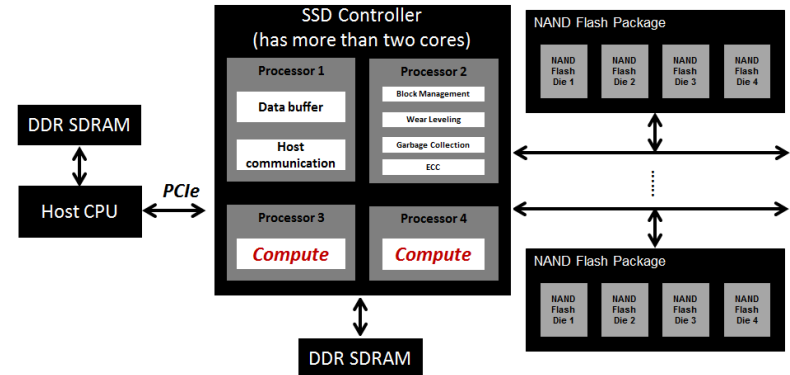
[kevin.gomez@seagate.com](mailto:kevin.gomez@seagate.com)

# Architecture Block Diagrams

Baseline – SSD for Data, Compute in Host

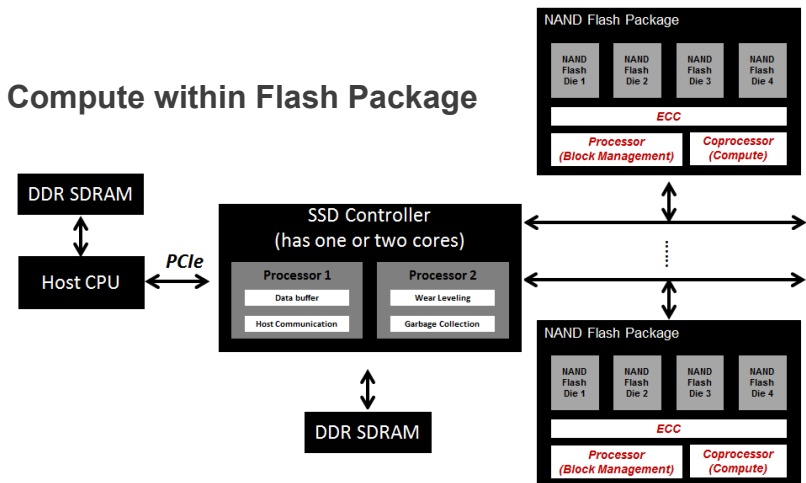


“Active Flash” – Compute in SSD Controller Processor



Added Compute functionality may be power islanded and enabled through firmware to make an SCE indistinguishable from a standard SSD - sharing the same production flow and economy of scale

Compute within Flash Package





# Architecture Simulation Parameters

CPU Power: use ITRS HP technology to evaluate dynamic and leakage power.

Number of Gates: 200M/core

Frequency: 2GHz.

Dynamic Power (per core): 5.04W

Leakage Power (per core): 0.340W

SSD Controller Power: use ITRS LOP technology to evaluate dynamic and leakage power.

Number of Gates: 20 millions per core (Assumption: 10% of the CPU).

Frequency: 1GHz.

Dynamic Power (per core): 0.156W.

Leakage Power (per core): 1.34mW.

Channel Processor Power: use ITRS LOP technology to evaluate dynamic and leakage power.

Number of Gates: 1K, 10K, 100K, 1M.

Frequency: 400MHz.

Dynamic Power (per core): 3.12uW, 31.2uW, 312uW, 3.12mW.

Leakage Power (per core): 67nW, 670nW, 6.7uW, 67uW.

DDR SDRAM: use parameters from MICRON.

Dynamic Power (per 2GB): 438.3mW.

Leakage Power (per 2GB): 88.1mW.

NAND Flash: use parameters from MICRON.

Dynamic Power (per die): 0.04W.

Leakage Power (per die): 0.003W.

Host Interface: PCIe.

Dynamic Power (per GB): 37.5mW.

Leakage Power (per GB): 0.mW

# Baseline Face Recognition

|   | 1-Core | 2-Core | 4-Core | 8-Core | 16-Core |
|---|--------|--------|--------|--------|---------|
| <b>Average Processing Time of Facial Recognition Algorithm (ms)</b> |        |        |        |        |         |
| CPI = 100   | 52.7   | 26.4   | 13.3   | 6.80   | 3.50    |
| CPI = 10  | 5.50   | 2.90   | 3.40   | 3.00   | 2.90    |
| CPI = 1   | 3.00   | 2.80   | 2.80   | 2.70   | 2.70    |
| CPI = 0.1   | 2.70   | 2.70   | 2.70   | 2.70   | 2.70    |
| <b>Average Power of Facial Recognition Algorithm (W)</b>            |        |        |        |        |         |
| CPI = 100   | 5.58   | 10.93  | 21.5   | 42.1   | 81.5    |
| CPI = 10  | 5.73   | 10.86  | 9.97   | 12.32  | 15.64   |
| CPI = 1   | 2.12   | 2.54   | 3.27   | 4.67   | 7.44    |
| CPI = 0.1   | 1.39   | 1.74   | 2.43   | 3.81   | 6.56    |
| <b>Average Energy of Facial Recognition Algorithm (mJ)</b>          |        |        |        |        |         |
| CPI = 100   | 294    | 289    | 286    | 287    | 286     |
| CPI = 10  | 31.5   | 31.5   | 33.9   | 36.9   | 45.4    |
| CPI = 1   | 6.36   | 7.13   | 9.17   | 12.6   | 20.1    |
| CPI = 0.1   | 3.76   | 4.71   | 6.57   | 10.3   | 17.7    |

Core = Host CPU Cores  
CPI = clock cycles per instruction of single core in CPU

# Active Flash Face Recognition

|   | 1-Core | 2-Core | 4-Core | 8-Core | 16-Core |
|---|--------|--------|--------|--------|---------|
| <b>Average Processing Time of Facial Recognition Algorithm (ms)</b> |        |        |        |        |         |
| CPI = 100   | 52.6   | 26.4   | 13.3   | 6.70   | 3.40    |
| CPI = 10  | 5.40   | 2.80   | 1.50   | 0.800  | 0.500   |
| CPI = 1   | 0.700  | 0.400  | 0.300  | 0.300  | 0.300   |
| <b>Average Power of Facial Recognition Algorithm (W)</b>            |        |        |        |        |         |
| CPI = 100   | 0.699  | 0.858  | 1.17   | 1.79   | 2.98    |
| CPI = 10  | 0.716  | 0.881  | 1.18   | 1.70   | 2.48    |
| CPI = 1   | 0.839  | 1.02   | 1.15   | 1.16   | 1.17    |
| <b>Average Energy of Facial Recognition Algorithm (mJ)</b>          |        |        |        |        |         |
| CPI = 100   | 36.8   | 22.6   | 15.6   | 12.0   | 10.1    |
| CPI = 10  | 3.86   | 2.47   | 1.78   | 1.36   | 1.24    |
| CPI = 1   | 0.587  | 0.410  | 0.345  | 0.347  | 0.351   |

Core = SSD Controller Cores  
CPI = clock cycles per  
instruction of single core in  
SSD controller

Exploiting Free Silicon for Energy-Efficient Computing Directly in NAND Flash-based Solid-State Storage Systems, High Performance Extreme Computing 2013, Li et al

# In-Flash-Package Face Recognition

| Channels   | 4     | 8     | 16    | 32     |
|--|-------|-------|-------|--------|
| <b>Average Processing Time (ms)</b>                        |       |       |       |        |
| Time   | 0.300 | 0.200 | 0.100 | 0.0500 |
| <b>Average Power of Facial Recognition Algorithm (W)</b>   |       |       |       |        |
| Gates = 1K   | 0.887 | 1.23  | 1.87  | 2.98   |
| Gates = 10K  | 0.887 | 1.23  | 1.87  | 2.98   |
| Gates = 100K   | 0.888 | 1.23  | 1.88  | 2.99   |
| Gates = 1M   | 0.899 | 1.26  | 1.92  | 3.06   |
| <b>Average Energy of Facial Recognition Algorithm (mJ)</b> |       |       |       |        |
| Gates = 1K   | 0.266 | 0.246 | 0.187 | 0.149  |
| Gates = 10K  | 0.266 | 0.246 | 0.187 | 0.149  |
| Gates = 100K   | 0.266 | 0.247 | 0.188 | 0.149  |
| Gates = 1M   | 0.270 | 0.251 | 0.192 | 0.153  |

# References

- [1] “Hitting the memory wall: implications of the obvious”, WA Wulf, SA McKee - ACM SIGARCH computer architecture news, 1995
- [2] “Reflections on the memory wall”, SA McKee - Proceedings of the 1st conference on Computing, 2004
- [3] “Missing the memory wall: The case for processor/memory integration”, A Nowatzky, F Pong, A Saulsbury - Architecture, 1996
- [4] “Computing performance: Game over or next level?”, SH Fuller, LI Millett - Computer, 2011
- [5] “Platform 2015: Intel processor and platform evolution for the next decade”, S Borkar, P Dubey, K Kahn, D Kuck, H Mulder
- [6] “Dark silicon and the end of multicore scaling”, H Esmaeilzadeh, E Blem, RS Amant... (ISCA), 2011
- [7] “GPUs and the future of parallel computing”, SW Keckler, WJ Dally, B Khailany, M Garland - Micro, 2011
- [8] “The GPU computing era”, J Nickolls, WJ Dally - Micro, IEEE, 2010
- [9] “Architecture at the End of Moore”, S Kaxiras - 2013 – Springer
- [10] “The Shift to Cloud Computing: Forget the Technology, It’s About Economics”, Jim Cooke 2010
- [11] “An Energy-Efficient Processor Architecture for Embedded Systems”, Balfour et. al. 2008
- [12] “Trends in Computation, Communication and Storage and the Consequences for Data-intensive Science”, Oliveira, S.F., 2012